



Andrieu, C., Yildirim, S., Doucet, A., & Chopin, N. (2020). Metropolis-Hastings with Averaged Acceptance Ratios. *arXiv*.  
<https://arxiv.org/abs/2101.01253>

Early version, also known as pre-print

[Link to publication record in Explore Bristol Research](#)  
PDF-document

## University of Bristol - Explore Bristol Research

### General rights

This document is made available in accordance with publisher policies. Please cite only the published version using the reference above. Full terms of use are available:  
<http://www.bristol.ac.uk/red/research-policy/pure/user-guides/ebr-terms/>

# Metropolis-Hastings with Averaged Acceptance Ratios

Christophe Andrieu<sup>\*</sup>, Sinan Yıldırım<sup>+</sup>, Arnaud Doucet<sup>†</sup>, and Nicolas Chopin<sup>◆</sup>

January 6, 2021

<sup>\*</sup>School of Mathematics, University of Bristol, U.K.

<sup>†</sup>Department of Statistics, University of Oxford, U.K.

<sup>+</sup>Faculty of Engineering and Natural Sciences, Sabancı University, Turkey.

<sup>◆</sup>ENSAE, France.

## Abstract

Markov chain Monte Carlo (MCMC) methods to sample from a probability distribution  $\pi$  defined on a space  $(\Theta, \mathcal{T})$  consist of the simulation of realisations of Markov chains  $\{\theta_n, n \geq 1\}$  of invariant distribution  $\pi$  and such that the distribution of  $\theta_i$  converges to  $\pi$  as  $i \rightarrow \infty$ . In practice one is typically interested in the computation of expectations of functions, say  $f$ , with respect to  $\pi$  and it is also required that averages  $M^{-1} \sum_{n=1}^M f(\theta_n)$  converge to the expectation of interest. The iterative nature of MCMC makes it difficult to develop generic methods to take advantage of parallel computing environments when interested in reducing time to convergence. While numerous approaches have been proposed to reduce the variance of ergodic averages, including averaging over independent realisations of  $\{\theta_n, n \geq 1\}$  simulated on several computers, techniques to reduce the “burn-in” of MCMC are scarce. In this paper we explore a simple and generic approach to improve convergence to equilibrium of existing algorithms which rely on the Metropolis-Hastings (MH) update, the main building block of MCMC. The main idea is to use averages of the acceptance ratio w.r.t. multiple realisations of random variables involved, while preserving  $\pi$  as invariant distribution. The methodology requires limited change to existing code, is naturally suited to parallel computing and is shown on our examples to provide substantial performance improvements both in terms of convergence to equilibrium and variance of ergodic averages. In some scenarios gains are observed even on a serial machine.

Keywords: Doubly intractable distributions; Intractable likelihood; Markov chain Monte Carlo; Pseudo-marginal Metropolis-Hastings; Reversible jump Monte Carlo; Sequential Monte Carlo; State-space models; Particle MCMC.

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	Estimating the target density . . . . .	4
1.2	Estimating the acceptance ratio . . . . .	4
1.3	Contribution . . . . .	5
<b>2</b>	<b>Using averaged acceptance ratio estimators</b>	<b>6</b>
2.1	A general perspective on MH based algorithms . . . . .	6
2.2	Motivation: an idealised algorithm . . . . .	7
2.3	MH with Averaged Acceptance Ratio . . . . .	9
2.4	Example: exchange algorithm and some analysis . . . . .	11
2.5	Example: improving transdimensional samplers . . . . .	13
<b>3</b>	<b>An efficient alternative to pseudo-marginal algorithms</b>	<b>17</b>
3.1	A novel consistent pseudo-marginal estimator . . . . .	17
3.2	Examples . . . . .	20
<b>4</b>	<b>State-space models: SMC and cSMC within MHAAR</b>	<b>24</b>
4.1	State-space models and cSMC . . . . .	25
4.2	MHAAR with cSMC for state-space models . . . . .	26
4.2.1	Unbiased estimator of the acceptance ratio using particles of cSMC . . . . .	26
4.2.2	MHAAR-RB for SSM . . . . .	27
4.2.3	Reduced computational cost via subsampling . . . . .	31
<b>5</b>	<b>Discussion</b>	<b>31</b>
<b>6</b>	<b>Acknowledgements</b>	<b>32</b>
<b>A</b>	<b>Proofs for the theorems in Section 2</b>	<b>35</b>
A.1	Acceptance ratio of Algorithm 1 . . . . .	35
A.2	Proof of Theorem 2 . . . . .	36
<b>B</b>	<b>Proofs for Section 3</b>	<b>37</b>
B.1	Acceptance ratio of Algorithm 3 . . . . .	37
B.2	Delayed rejection step for Algorithm 3 . . . . .	39
<b>C</b>	<b>Auxiliary results and proofs Section 4</b>	<b>42</b>
C.1	Unbiasedness for the acceptance ratio estimator of Algorithm 5 . . . . .	45
C.2	Proofs for Algorithm 5 . . . . .	46
C.2.1	Acceptance ratio of Algorithm 5 . . . . .	46
C.2.2	Delayed rejection step for Algorithm 5 . . . . .	48
C.3	The subsampled version of MHAAR-RB for SSM . . . . .	50
C.3.1	Reversibility of Algorithm 7 . . . . .	50
C.3.2	Refreshing the latent variable in Algorithm 7 . . . . .	53

# 1 Introduction

Suppose we wish to sample from a given probability distribution  $\pi$  on some measurable space  $(\Theta, \mathcal{T})$ . When it is impossible or too difficult to generate perfect samples from  $\pi$ , one practical resource is to use a Markov chain Monte Carlo (MCMC) algorithm generating an ergodic Markov chain  $\{\theta_n, n \geq 0\}$  whose invariant distribution is  $\pi$ . Among MCMC methods, the Metropolis–Hastings (MH) algorithm plays a central rôle. The MH update proceeds as follows: given  $\theta_n = \theta$  and a Markov transition kernel  $q(\theta, \cdot)$  on  $(\Theta, \mathcal{T})$ , we propose  $\vartheta \sim q(\theta, \cdot)$  and set  $\theta_{n+1} = \vartheta$  with probability  $\alpha(\theta, \vartheta) := \min\{1, r(\theta, \vartheta)\}$ , where

$$r(\theta, \vartheta) := \frac{\pi(d\vartheta)q(\vartheta, d\theta)}{\pi(d\theta)q(\theta, d\vartheta)} \quad (1)$$

for  $(\theta, \vartheta) \in \mathbf{S} \subset \Theta^2$  (see Tierney [1998] for a definition of  $\mathbf{S}$ ) is a well defined Radon–Nikodym derivative, and  $r(\theta, \vartheta) = 0$  otherwise. When the proposed value is rejected, we set  $\theta_{n+1} = \theta$ . We will refer to  $r(\theta, \vartheta)$  as the acceptance ratio. The transition kernel of the Markov chain  $\{\theta_n, n \geq 0\}$  generated with the MH algorithm with proposal kernel  $q(\cdot, \cdot)$  is

$$P(\theta, A) = \int_A \alpha(\theta, \vartheta)q(\theta, d\vartheta) + \rho(\theta)\mathbb{I}\{\theta \in A\}, \quad (\theta, A) \in \Theta \times \mathcal{T}, \quad (2)$$

where  $\rho(\theta)$  is the rejection probability such that  $P(\theta, \Theta) = 1$  and  $\mathbb{I}\{\cdot \in A\}$  is the indicator function for set  $A$ . Expectations of functions, say  $f$ , with respect to  $\pi$  can be estimated with  $S_M := M^{-1} \sum_{n=1}^M f(\theta_n)$  for  $M \in \mathbb{N}$ , which is consistent under mild assumptions.

Being able to evaluate the acceptance ratio  $r(\theta, \vartheta)$  is therefore central to implementing the MH algorithm in practice. Recently, there has been much interest in expanding the scope of the MH algorithm to situations where this acceptance ratio is intractable, that is, impossible or very expensive to compute. A canonical example of intractability is when  $\pi$  can be written as the marginal of a given joint probability distribution for  $\theta$  and some latent variable  $z$ . A classical way of addressing this problem consists of running an MCMC algorithm targeting the joint distribution, which may however become very inefficient in situations where the size of the latent variable is high—this is for example the case for general state-space models. In what follows, we will briefly review generic ways of tackling this problem. To that purpose we will use the following simple running example to illustrate various methods. This example has the advantage that its setup is relatively simple and of clear practical relevance. We postpone developments for much more complicated scenarios to Sections 2, 3, and 4.

**Example 1 (Inference with doubly intractable models).** In this scenario the likelihood function of the unknown parameter  $\theta \in \Theta$  for the dataset  $y \in \mathbf{Y}$ ,  $\ell_\theta(y)$ , is only known up to a normalising constant, that is  $\ell_\theta(y) = g_\theta(y)/C_\theta$ , where  $C_\theta$  is unknown, while  $g_\theta(y)$  can be evaluated pointwise for any value of  $\theta \in \Theta$ . In a Bayesian framework, for a prior density  $\eta(\theta)$ , we are interested in the posterior density  $\pi(\theta)$ , given by  $\pi(\theta) \propto \eta(\theta)\ell_\theta(y)$ . The acceptance ratio of the MH algorithm associated to a proposal density  $q(\theta, \vartheta)$  is

$$r(\theta, \vartheta) = \frac{q(\vartheta, \theta) \eta(\vartheta) g_\vartheta(y) C_\theta}{q(\theta, \vartheta) \eta(\theta) g_\theta(y) C_\vartheta}, \quad (3)$$

which cannot be calculated because of the unknown ratio  $C_\theta/C_\vartheta$ . While the likelihood function may be intractable, sampling artificial datasets  $u \sim \ell_\theta(y_*)dy_*$  may be possible

for any  $\theta \in \Theta$ , and sometimes computationally cheap. We will describe two known approaches which exploit and expand this property in order to design Markov kernels preserving  $\pi(\theta)$  as invariant density.

## 1.1 Estimating the target density

Assume for simplicity of exposition that  $\pi$  has a probability density with respect to some  $\sigma$ -finite measure. We will abuse notation slightly by using  $\pi$  for both the probability distribution and its density. A simple method to tackle intractability which has recently attracted interest consists of replacing the value of  $\pi(\theta)$  with a non-negative random estimator  $\hat{\pi}(\theta)$  whenever it is required in the implementation of the MH algorithm above. If there exists a constant  $C > 0$  such that  $\mathbb{E}[\hat{\pi}(\theta)] = C\pi(\theta)$  for all  $\theta \in \Theta$ , a property we refer to abusively as unbiasedness, this strategy turns out to lead to exact algorithms, that is sampling from  $\pi$  is guaranteed at equilibrium under very mild assumptions on  $\hat{\pi}(\theta)$ . This approach leads to so called pseudo-marginal algorithms [Beaumont, 2003, Andrieu and Roberts, 2009]. In what follows, for  $a, b \in \mathbb{R}$  we let  $\llbracket a, b \rrbracket := [a, b] \cap \mathbb{Z}$  and use the specialised notation  $\llbracket a \rrbracket := \llbracket 1, a \rrbracket$ .

**Example 2 (Example 1, ctd).** Let  $h : \mathcal{Y} \rightarrow [0, \infty)$  be an integrable non-negative function of integral equal to 1. For a given  $\theta$ , an unbiased estimate of  $\pi(\theta)$  can be obtained via importance sampling whenever the support of  $g_\theta$  includes that of  $h$ :

$$\hat{\pi}^N(\theta) \propto \eta(\theta)g_\theta(y) \left\{ \frac{1}{N} \sum_{i=1}^N \frac{h(u^{(i)})}{g_\theta(u^{(i)})} \right\}, \quad u^{(i)} \stackrel{\text{iid}}{\sim} \ell_\theta(y_*)dy_*, \quad i \in \llbracket N \rrbracket, \quad (4)$$

since the normalised sum is an unbiased estimator of  $1/C_\theta$ . The auxiliary variable method of Møller et al. [2006] corresponds to  $N = 1$ . An interesting feature of this approach is that  $N$  is a free parameter of the algorithm which reduces the variability of this estimator. It is shown in Andrieu and Vihola [2016] that increasing  $N$  always reduces the asymptotic variance of averages using this chain and will in most cases of interest improve convergence to equilibrium. This is particularly interesting in a parallel computing environment but, as we shall see, can prove of interest on serial machines.

## 1.2 Estimating the acceptance ratio

One can in fact push the idea of replacing algebraic expressions with estimators further. Instead of approximating the numerator and denominator of the acceptance ratio  $r(\theta, \vartheta)$  independently, it is indeed possible to use directly estimators of the acceptance ratio  $r(\theta, \vartheta)$  and still obtain algorithms guaranteed to sample from  $\pi$  at equilibrium. An interesting feature of these algorithms is that we estimate the ratio  $r(\theta, \vartheta)$  afresh whenever it is required. On the contrary, in algorithms using unbiased estimates of the target density, the estimate  $\hat{\pi}(\theta)/C$  is used in the acceptance ratio until a transition is accepted. As a consequence whenever  $\hat{\pi}(\theta)/C$  significantly overestimates  $\pi(\theta)$  the algorithm spends a long period of time stuck in a particular state, resulting in poor performance. In the following continuation of Example 1, we present a particular case of estimating the acceptance ratio, proposed by Murray et al. [2006].

**Example 3 (Example 1, ctd).** The exchange algorithm of [Murray et al. \[2006\]](#) is motivated by the realisation that while for  $u \sim \ell_\vartheta(y_*)dy_*$  and  $h(u)/g_\vartheta(u)$  is an unbiased estimator of  $1/C_\vartheta$ , the particular choice  $h(u) = g_\theta(u)$  leads to an unbiased estimator  $g_\theta(u)/g_\vartheta(u)$  of  $C_\theta/C_\vartheta$  required in (3). This suggests the following MH type update. Given  $\theta \in \Theta$ , sample  $\vartheta \sim q(\theta, \cdot)$ , then  $u \sim \ell_\vartheta(y_*)dy_*$  and use the acceptance ratio

$$r_u(\theta, \vartheta) = \frac{q(\vartheta, \theta) \eta(\vartheta) g_\vartheta(y) g_\theta(u)}{q(\theta, \vartheta) \eta(\theta) g_\theta(y) g_\vartheta(u)}, \quad (5)$$

which is an unbiased estimator of the acceptance ratio in (3). Remarkably this algorithm admits  $\pi$  as an invariant distribution and hence, under additional mild assumptions, is guaranteed to produce samples asymptotically distributed according to  $\pi$ .

### 1.3 Contribution

As we shall see numerous MH algorithms of interest to sample from  $\pi$  have a tractable acceptance ratio of the form  $r_u(\theta, \vartheta)$  where  $u$  is sampled afresh at each iteration, as is the case in Example 3. Such sampling induces variability of the acceptance ratio which, as we shall see, is undesirable and a natural question is whether this can be alleviated by averaging multiple realisations of some of the variables involved. More specifically, given multiple realisations  $r_{u(i)}(\theta, \vartheta)$ ,  $i \in \llbracket N \rrbracket$  is it possible to design an algorithm leaving  $\pi$  invariant and of superior performance? While the naïve approach consisting of using  $N^{-1} \sum_{i=1}^N r_{u(i)}(\theta, \vartheta)$  in place of  $r_u(\theta, \vartheta)$  is not valid, in that  $\pi$  is not guaranteed to be an invariant distribution anymore, we show that a solution alternating between the use of this average and its inverse leads to a correct algorithm. These algorithms naturally lend themselves to parallel computations as independent ratio estimators can be computed in parallel at each iteration [Lee et al. \[2010\]](#), [Suchard et al. \[2010\]](#). Provided access to a parallel machine is available and the cost of computing  $r_u(\theta, \vartheta)$  dominates communication cost, which is the case in challenging applications, we show that this approach can reduce the burn-in-period, sometimes substantially—in fact the higher the variability of  $r_u(\theta, \vartheta)$  for  $\theta, \vartheta \in \Theta$  the more substantial the gains are. As a by-product the induced rapid mixing also leads to reduced asymptotic variance of ergodic averages, even when implemented on a serial machine in some scenarios. Generic methods to reduce burn-in and utilise parallel architectures are scarce [\[Sohn, 1995\]](#), in contrast with variance reduction techniques for which better embarrassingly parallel solutions [\[Sherlock et al., 2017, Bornn et al., 2017\]](#) and/or post-processing methods are available [\[Delmas and Jourdain, 2009, Dellaportas and Kontoyiannis, 2012\]](#). An interesting practical point is that the approach we advocate requires only limited adaptation of the specific, and often intricate, code for an existing algorithm beyond the generic management of the parallel environment. Note however that the actual implementation of our algorithms on a parallel computer is beyond of the present manuscript which focuses primarily on developing sound methodology and provide initial evaluation of expected performance.

In Sections 2 we introduce the MHAAR methodology in full generality, providing some theoretical analysis supporting their correctness and claimed efficiency while we illustrate its interest in the context of reversible jump MCMC algorithms. In Section 3 we specialise MHAAR to latent variable models and present an alternative to pseudo-marginal algorithms [Beaumont \[2003\]](#), [Andrieu and Roberts \[2009\]](#) which is shown to

have far superior performance properties, even on a serial machine. In Section 4, we show how MHAAR can be advantageous in the context of inference in state-space models when it is utilised in combination with sequential Monte Carlo (SMC) algorithms. In particular, we expand the scope of particle MCMC algorithms [Andrieu et al., 2010] and show novel ways of using multiple or all possible paths obtainable from a conditional SMC (cSMC) run to estimate the marginal acceptance ratio. of MHAAR. We again assess gain performance numerically, demonstrating the interest of the approach. The proofs of the validity of our algorithms as well as additional discussion on the generalisation of the methods can found in the Appendices.

## 2 Using averaged acceptance ratio estimators

### 2.1 A general perspective on MH based algorithms

Before describing our novel algorithms we briefly outline a framework, fully developed in Andrieu et al. [2020], which allows for a systematic and concise presentation of complex MH updates. In particular the presentation adopted makes validating, that is establishing reversibility with respect to the distribution of interest, fairly direct and is helpful to establish the expression for the acceptance ratio involved in the update.

The key idea here is that in order to describe and validate a MH update it is sufficient to identify all the random variables  $\xi$  involved in the update before the accept/reject step, their distribution, the mapping  $\varphi$  used to determine the next state of the Markov chain from  $\varphi(\xi)$  and check that it satisfies  $\varphi \circ \varphi = \text{Id}$ . Consider for example the standard update given at the beginning of Section 1: here the variables involved are  $\xi := (\theta, \vartheta) \in \Theta^2$ , their distribution before the accept/reject step is  $\dot{\pi}(d\xi) = \pi(d\theta)q(d\vartheta)$ , and the involution used to determine the next state is  $\varphi(\theta, \vartheta) := (\vartheta, \theta)$  for  $\theta, \vartheta \in \Theta^2$ , leading to the familiar acceptance ratio (1). The popular random walk Metropolis algorithm corresponds to the choices  $\xi = (\theta, \zeta) \in \Theta^2$ , where  $\zeta \in \Theta$  is the increment used to perturb  $\theta$ ,  $\dot{\pi}(d\xi) = \pi(d\theta)q(d\zeta)$  and  $\varphi(\theta, \zeta) = (\theta + \zeta, -\zeta)$ . In the situation where  $\Theta = \mathbb{R}^d$ ,  $\pi$  and  $q$  admit densities with respect to the Lebesgue measure (also denoted  $\pi$  and  $q$  and assumed to be strictly positive for simplicity) and  $q$  is symmetric, the resulting acceptance ratio is of form

$$\frac{\pi(\theta + \zeta)q(-\zeta)}{\pi(\theta)q(\zeta)} = \frac{\pi(\theta + \zeta)}{\pi(\theta)},$$

where the numerator is the density resulting from the change of variable  $\varphi(\theta, \zeta) = \varphi^{-1}(\theta, \zeta) = (\theta + \zeta, -\zeta)$ , of Jacobian 1. This can be generalised as follows. Let  $\dot{\pi}$  be a probability distribution on some measurable space  $(\mathbf{X}, \mathcal{X})$  and let  $\varphi : \mathbf{X} \rightarrow \mathbf{X}$  be a measurable mapping, we define the push forward distribution  $\dot{\pi}^\varphi$  to be the probability distribution of  $\varphi(\xi)$  when  $\xi \sim \dot{\pi}$ , that is such that for any measurable  $A \in \mathcal{X}$ ,  $\dot{\pi}^\varphi(A) := \dot{\pi}(\varphi^{-1}(A))$ . Assume further that  $\dot{\pi}$  has marginal  $\pi$ , say  $\dot{\pi}(d\xi) = \pi(d\xi_0)\dot{\pi}(d\xi_1 | \xi_0)$ , and that  $\varphi : \mathbf{X} \rightarrow \mathbf{X}$  is an involution. Then the following update is a valid MH update, that is ignoring the second components  $\xi_1$  and  $\xi'_1$  it is reversible with respect to  $\pi$  and hence leaves this distribution invariant:

1. given  $\xi_0$  sample  $\xi_1 \sim \dot{\pi}(\cdot | \xi_0)$ ,

2. compute

$$\alpha(\xi) := \min \{1, \mathring{r}(\xi)\} \text{ with } \mathring{r}(\xi) := \frac{\mathring{\pi}^\varphi(d\xi)}{\mathring{\pi}(d\xi)}, \quad (6)$$

3. with probability  $\alpha(\xi)$  return  $\xi' = \varphi(\xi)$ , otherwise return  $\xi' = \xi$ .

The quantity  $r(\xi)$  is a so-called Radon-Nikodym derivative, guaranteed to exist under very mild assumptions. In this manuscript  $\mathring{\pi}$  will always be assumed to have a known density with respect to a product of counting and Lebesgue measures and  $r(\xi)$  will be either zero whenever either densities of  $\xi' = \varphi(\xi)$  or  $\xi$  is zero, or the ratio of these densities otherwise. The notation above allows us, for the moment, to avoid the distinction between discrete and real valued variables and the possible presence of a Jacobian. Naturally another practical requirement is that sampling from the “proposal distribution”  $\mathring{\pi}(\cdot \mid \xi_0)$  should be computationally tractable.

To summarize, in what follows we adopt the following systematic presentation of MH updates:

1. identify all the instrumental variables  $\xi_1$  and the distribution  $\mathring{\pi}(d\xi)$  involved in the parameter update,
2. identify the involution  $\varphi: \mathbf{X} \rightarrow \mathbf{X}$ ,
3. find an expression for  $\mathring{r}(\xi)$ .

Note that the above does not ensure convergence to equilibrium of the Markov chain, which is problem dependent.

The following property can be established and will be used on several occasions in the remainder of the manuscript, for  $\xi \in \mathbf{X}$

$$\mathring{r} \circ \varphi(\xi) = \begin{cases} 1/\mathring{r}(\xi) & \text{if } \mathring{r}(\xi) > 0 \\ 0 & \text{otherwise} \end{cases}. \quad (7)$$

In order to simplify presentation we will always assume that  $\mathring{r}(\xi) > 0$  for any  $\xi \in \mathbf{X}$ —the general scenario is a straightforward adaptation.

## 2.2 Motivation: an idealised algorithm

Consider the generic algorithm given in the previous subsection. Our primary aim here is to show that it is possible to improve performance of this algorithm by using a modification where the acceptance ratio  $\mathring{r}(\xi)$  in (6) is integrated with respect to a subset of the proposed variables  $\xi_1$ . In the case of Example 1-3, we have  $\xi_1 = (\vartheta, u) \in \Theta \times \mathbf{Y}$  and marginalisation with respect to  $u$ , that is the simulated artificial datasets, is sought. The motivation for this is that removing dependence of  $\mathring{r}(\xi)$  on  $u$  removes variability and will result in a better expected acceptance rate and, in the spirit of [Andrieu and Vihola, 2016] lead to algorithms of improved performance. The algorithm is not implementable in general but captures in a simple setup the main idea we develop further in this paper. Indeed MHAAR algorithms are exact numerical approximations of this idealised algorithm, in that they preserve the desired distribution invariant, and the latter algorithm



can be thought of as a ‘lower bound’ on what the approximations can achieve in terms of performance.

Motivated by applications, we consider the scenario where the target distribution of interest  $\pi(d\theta)$  is not tractable, but arises from a tractable latent variable model  $\pi(d(\theta, z))$  defined on some space  $(\Theta \times \mathbf{Z}, \mathcal{T} \otimes \mathcal{Z})$ . As a result the target distribution of interest is now  $\pi(d(\theta, z))$  and Example 1 can be recovered by simply ignoring  $z$ . We first describe a standard instance of the MH update to sample from this target. Let  $(\mathbf{U}, \mathcal{U})$  be some probability space, and  $\phi_{\theta, \vartheta} : \mathbf{Z} \times \mathbf{U} \mapsto \mathbf{Z} \times \mathbf{U}$  for all  $\theta, \vartheta \in \Theta^2$  be invertible mappings such that  $\phi_{\theta, \vartheta} = \phi_{\vartheta, \theta}^{-1}$ . Using the framework of the previous section we consider the set of variables  $\xi := (\theta, \vartheta, z, u) \in \Theta^2 \times \mathbf{Z} \times \mathbf{U}$ , an involution of the type

$$\varphi(\theta, \vartheta, z, u) := (\vartheta, \theta, \phi_{\theta, \vartheta}(z, u)), \quad (8)$$

and the probability distribution

$$\tilde{\pi}(d(\theta, \vartheta, z, u)) := \pi(d(\theta, z))q(\theta, d\vartheta)Q_{\theta, \vartheta, z}(du), \quad (9)$$

for a family of probability distributions  $Q_{\theta, \vartheta, z}(du)$  defined on the probability space  $(\mathbf{U}, \mathcal{U})$  and  $q(\theta, \cdot)$  as in Section 1. Here the nature of  $u$  is problem dependent, guided by the choice of involution  $\varphi$  and tractability of acceptance ratios of the type (6). This generality allows us to cover scenarios where the latent variable  $z$  is updated thanks to a mapping from  $z, u$  to  $z', u'$  by  $\phi_{\theta, \vartheta}(\cdot)$ . For example, again ignoring the latent variable  $z$  from the notation and letting  $u' = u \in \mathbf{U} = \mathbf{Y}$  corresponds to the exchange algorithm of Example 1-3. We now introduce an improved MH update which uses the integrated acceptance ratio

$$\hat{r}(\theta, \vartheta, z) := \int \hat{r}(\theta, \vartheta, z, u)Q_{\theta, \vartheta, z}(du), \quad (10)$$

assumed to be tractable for the moment. We note that in Example 1-3 the choice  $Q_{\theta, \vartheta, z}(du) = \ell_{\vartheta}(u)du$ , where we keep  $z$  for notational compatibility but recall that  $z$  is not needed for this example, the integrated acceptance ratio simplifies to (3). A solution around intractability is the topic of the next section. The following update can be shown to be  $\tilde{\pi}$ -reversible.

1. sample  $\vartheta \sim q(\theta, \cdot)$  and  $c \sim \text{Unif}\{1, 2\}$
2. sample

$$u \sim \begin{cases} Q_{\theta, \vartheta, z}(du) \hat{r}(\theta, \vartheta, z, u) / \hat{r}(\theta, \vartheta, z) & \text{if } c = 1 \\ Q_{\theta, \vartheta, z}(du) & \text{if } c = 2 \end{cases} \quad (11)$$

and form  $\xi := (\theta, \vartheta, z, u, c)$ ,

3. with  $\varphi$  as above, compute  $\xi' = \varphi^*(\xi) := (\varphi(\theta, \vartheta, z, u), 3 - c) = (\vartheta, \theta, z', u', 3 - c)$ ,
4. return  $\xi'$  with probability  $\alpha(\theta, \vartheta, z, u, c) = \min\{1, \hat{r}(\theta, \vartheta, z, u, c)\}$  where

$$\hat{r}(\theta, \vartheta, z, u, c) = \begin{cases} \hat{r}(\theta, \vartheta, z) & \text{if } c = 1 \\ 1 / \hat{r}(\vartheta, \theta, z') & \text{if } c = 2 \end{cases},$$

otherwise return  $\xi = (\theta, \vartheta, z, u, c)$ .

The essential idea here is that alternating between the use of two appropriately chosen sampling schemes for  $u$ , the acceptance probability depends on the integrated acceptance ratio only. What's more in the case where  $c = 1$  we see that the proposal distribution for  $u$  is biased towards values leading to high acceptance ratios for the algorithm defined by (8) and (9), in the spirit of [Cainey \[2013, Chapter 4\]](#) and [Zanella \[2020\]](#) where the proposal distribution is weighted by a function of the target density  $\pi$ . We now briefly outline why the acceptance ratio appears to be integrated:

$$\hat{\pi}^*(d(\theta, \vartheta, z, u, c)) = \begin{cases} \hat{\pi}(d(\theta, \vartheta, z, u)) \hat{r}(\theta, \vartheta, z, u) / \hat{r}(\theta, \vartheta, z)^{\frac{1}{2}} & \text{if } c = 1 \\ \hat{\pi}(d(\theta, \vartheta, z, u))^{\frac{1}{2}} & \text{if } c = 2 \end{cases},$$

using (6) we see that the acceptance ratio is of the form claimed, as for  $c = 1$

$$\frac{(\hat{\pi}^*)^\varphi(d(\theta, \vartheta, z, u, 2))}{\hat{\pi}^*(d(\theta, \vartheta, z, u, 1))} = \frac{\hat{\pi}^\varphi(d(\theta, \vartheta, z, u))}{\hat{\pi}(d(\theta, \vartheta, z, u))} \frac{\hat{r}(\theta, \vartheta, z)}{\hat{r}(\theta, \vartheta, z, u)} = \hat{r}(\theta, \vartheta, z),$$

which does not depend on  $u$ , and for  $c = 2$  we use (7), yielding  $\hat{r}(\theta, \vartheta, z, u, 2) = 1/\hat{r} \circ \varphi^*(\theta, \vartheta, z, u, 1) = 1/\hat{r}(\vartheta, \theta, z')$ , which depends on  $u$  through  $z'$ . The acceptance ratio for  $c = 2$  may seem disappointing, but it can be shown that reversibility implies

$$\int \alpha(\theta, \vartheta, z, u, 1) \hat{\pi}^*(d(\theta, \vartheta, z, u, 1)) = \int \alpha(\theta, \vartheta, z, u, 2) \hat{\pi}^*(d(\theta, \vartheta, z, u, 2)),$$

that is the expected acceptance probabilities when  $c = 1$  or  $c = 2$  are equal. Further application of Jensen's inequality to the concave function  $a \mapsto \min\{1, a\}$  establishes that

$$\frac{1}{2} \int \min\{1, \hat{r}(\theta, \vartheta, z, u)\} \hat{\pi}(d(\theta, \vartheta, z, u)) \leq \int \min\{1, \hat{r}(\theta, \vartheta, z)\} \hat{\pi}^*(d(\theta, \vartheta, z, u, 1)),$$

implying that for a given proposal mechanism  $q(\theta, d\vartheta)$ , the algorithm using the integrated acceptance ratio accepts more proposed transitions. We will see that this leads to improved performance.

## 2.3 MH with Averaged Acceptance Ratio

While valid theoretically, the algorithm of Subsection 2.2 is rarely implementable in practice since  $\hat{r}(\theta, \vartheta, z)$  is typically intractable and sampling  $u$  from (11) when  $c = 1$  potentially difficult. Instead we develop here a very closely related update relying on averages of

$$r_{u^{(i)}}(\theta, \vartheta, z) := \hat{r}(\theta, \vartheta, z, u^{(i)})$$

for, say,  $N > 1$  realisations  $\mathbf{u} := u^{(1:N)} = (u^{(1)}, \dots, u^{(N)}) \in \mathfrak{U} := \mathbf{U}^N$  of  $u$ , that is the accept/reject mechanism will now rely on

$$r_{\mathbf{u}}^N(\theta, \vartheta, z) := \frac{1}{N} \sum_{i=1}^N r_{u^{(i)}}(\theta, \vartheta, z). \quad (12)$$

The novel scheme, called MH with Averaged Acceptance Ratio (MHAAR), relies on the set of variables  $\xi := (\theta, \vartheta, z, \mathbf{u}, k, c) \in \Theta^2 \times \mathbf{Z} \times \mathfrak{U} \times \llbracket N \rrbracket \times \{1, 2\}$  and the joint distribution

$$\hat{\pi}(d\xi) := \pi(d(\theta, z)) \frac{1}{2} Q_c^N(\theta, z; d(\vartheta, \mathbf{u}, k)), \quad (13)$$

where the probability distributions  $Q_c^N(\theta, z; d(\vartheta, \mathbf{u}, k))$ ,  $c = 1, 2$ , are given by

$$Q_1^N(\theta, z; d(\vartheta, \mathbf{u}, k)) := q(\theta, d\vartheta) \prod_{i=1}^N Q_{\theta, \vartheta, z}(du^{(i)}) \frac{\frac{1}{N} r_{\mathbf{u}^{(k)}}(\theta, \vartheta, z)}{r_{\mathbf{u}}^N(\theta, \vartheta, z)},$$

$$Q_2^N(\theta, z; d(\vartheta, \mathbf{u}, k)) := q(\theta, d\vartheta) Q_{\theta, \vartheta, z}(du^{(k)}) \prod_{i=1, i \neq k}^N Q_{\vartheta, \theta, \phi_{\theta, \vartheta}^{[1]}(z, u^{(k)})}(du^{(i)}) \frac{1}{N},$$

where, with  $\phi_{\theta, \vartheta}$  as in Section 2.2, we have defined the functions  $\phi_{\theta, \vartheta}^{[1]}: \mathbf{Z} \times \mathbf{U} \rightarrow \mathbf{Z}$  and  $\phi_{\theta, \vartheta}^{[2]}: \mathbf{Z} \times \mathbf{U} \rightarrow \mathbf{U}$  such that  $\phi_{\theta, \vartheta} := (\phi_{\theta, \vartheta}^{[1]}, \phi_{\theta, \vartheta}^{[2]})$ . Here,  $r_{\mathbf{u}^{(k)}}(\theta, \vartheta, z)$  is the acceptance ratio corresponding to the joint distribution in (9) along with the involution in (8). As in the previous section, the role of  $\phi_{\theta, \vartheta}$  is to parametrise how a new value  $z'$  of  $z$  is proposed in an MH update using a variable  $u \in \mathbf{U}$  i.e.  $(z', u') = \phi_{\theta, \vartheta}(z, u)$ . A simple example corresponds to  $\mathbf{U} = \mathbf{Z}$  and the choice  $\phi_{\theta, \vartheta}(z, u) = (u, z)$ ; a more sophisticated example will be given in Section 2.5. A MHAAR update consists of the following steps. Given  $(\theta, z) \in \Theta \times \mathbf{Z}$ ,

1. sample  $c \sim \text{Unif}(\{1, 2\})$ ,  $(\vartheta, \mathbf{u}, k) \sim Q_c^N(\theta, z; \cdot)$  and form  $\xi := (\theta, \vartheta, z, \mathbf{u}, k, c)$ ,
2. compute

$$\xi' = \varphi(\xi) := (\vartheta, \theta, \phi_{\theta, \vartheta}^{[1]}(z, u^{(k)}), u^{(1:k-1)}, \phi_{\theta, \vartheta}^{[2]}(z, u^{(k)}), u^{(k+1:N)}, k, 3 - c), \quad (14)$$

3. return  $\xi'$  with probability  $\min\{1, \hat{r}(\xi)\}$  where with  $(\theta', \vartheta', z', \mathbf{u}', k', c') = \xi'$  and  $r_{\mathbf{u}}^N(\theta, \vartheta, z)$  given in (12),

$$\hat{r}(\xi) = \begin{cases} r_{\mathbf{u}}^N(\theta, \vartheta, z), & \text{for } c = 1 \\ 1/r_{\mathbf{u}'}(\vartheta, \theta, z'), & \text{for } c = 2 \end{cases}, \quad (15)$$

otherwise return  $\xi$ .

It is not difficult to check that the mapping  $\varphi$  is an involution, and Theorem 1 below establishes that  $\hat{r}(\xi)$  indeed simplifies to the desired form (15).

**Theorem 1.** *For the probability distribution  $\hat{\pi}$  and involution  $\varphi$  defined in (13) and (14) respectively the acceptance ratio  $\hat{r}(\xi)$  is as in (15).*

Details of the proof can be found in Appendix A.1 and pseudo-code is given in Algorithm 1—we will refer to the corresponding Markov kernel as  $\hat{P}^N$  for  $N \in \mathbb{N}_*$  with the simplification  $\hat{P}$  for  $N = 1$ . For  $w_1, w_2, \dots, w_m$  such that for  $m \in \mathbb{N}$ ,  $w_k \geq 0$  for  $k = 1, \dots, m$  and  $\sum_{k=1}^m w_k > 0$ , we define  $K \sim \mathcal{P}(w_1, \dots, w_m)$  to mean that  $\mathbb{P}(K = k) \propto w_k$ .

*Remark 1.* Note that when  $c = 2$  the simulation of  $k$  is in practice not required. Also, when  $c = 1$ , the acceptance probability does not depend on  $k$ , hence sampling of  $k$  is necessary only if the move is accepted, which can be exploited for faster implementation.

---

**Algorithm 1:** MHAAR for averaging PMR estimators

---

**Input:** Current sample  $(\theta, z)$   
**Output:** New sample

- 1 Sample  $\vartheta \sim q(\theta, \cdot)$  and  $c \sim \text{Unif}(\{1, 2\})$ .
- 2 **if**  $c = 1$  **then**
  - 3     **for**  $i = 1, \dots, N$  **do**
  - 4         Sample  $u^{(i)} \sim Q_{\theta, \vartheta, z}(\cdot)$
  - 5     Sample  $k \sim \mathcal{P}(r_{u^{(1)}}(\theta, \vartheta, z), \dots, r_{u^{(N)}}(\theta, \vartheta, z))$ , and set  $z' = \phi_{\theta, \vartheta}^{[1]}(z, u^{(k)})$ .
  - 6     Return  $(\vartheta, z')$  with probability  $\min\{1, r_u^N(\theta, \vartheta, z)\}$ , otherwise return  $(\theta, z)$ .
- 7 **else**
  - 8     Sample  $k \sim \text{Unif}(\llbracket N \rrbracket)$  and  $u^{(k)} \sim Q_{\theta, \vartheta, z}(\cdot)$ ,  $z' = \phi_{\theta, \vartheta}^{[1]}(z, u^{(k)})$ .
  - 9     **for**  $i = 1, \dots, N, i \neq k$ , **do**
  - 10         Sample  $u^{(i)} \sim Q_{\vartheta, \theta, z'}(\cdot)$ .
  - 11     Return  $(\vartheta, z')$  with probability  $\min\{1, 1/r_u^N(\vartheta, \theta, z')\}$ , otherwise return  $(\theta, z)$ .

---

*Remark 2.* In some scenarios, for given values  $(\theta, \vartheta) \in \Theta^2$  it may be preferable for computational efficiency to sample  $(\vartheta, u, k) \sim Q_c^N(\theta, z; \cdot)$  for  $c = 1$  rather than  $c = 2$ , or vice versa. This will be the case in Example 5. This is possible by changing the distribution of  $c$ : define a function  $\omega : \Theta^2 \times \{1, 2\} \rightarrow [0, 1]$  such that  $\omega(\theta, \vartheta, 1) + \omega(\theta, \vartheta, 2) = 1$ , therefore defining a probability distribution for  $c$ , for any  $(\theta, \vartheta) \in \Theta^2$ . The resulting averaged acceptance ratio is now

$$r_u^N(\theta, \vartheta, z) := \frac{\omega(\vartheta, \theta, 2)}{\omega(\theta, \vartheta, 1)} \frac{1}{N} \sum_{i=1}^N r_{u^{(i)}}(\theta, \vartheta, z). \quad (16)$$

*Remark 3.* We remark the link to some of the ideas developed in Zanella [2020], but also the differences in terms of what is being averaged and the fact that we are not constrained to finite discrete spaces.

We now turn to two illustrative examples.

## 2.4 Example: exchange algorithm and some analysis

The exchange algorithm [Murray et al., 2006] in Example 3 lends itself to acceptance ratio averaging and can serve to illustrate precisely the gains one may expect from the approach. Here the model does not involve the auxiliary variable  $z$ ,  $U = Y$ ,  $Q_{\theta, \vartheta}(\cdot)$  corresponds to  $\ell_{\vartheta}(\cdot)$  and  $\phi$  is the identity function on  $Y$ . We can therefore apply the MHAAR approach described in Algorithm 1. The algorithm takes the following form. Sample  $\vartheta \sim q(\theta, \cdot)$ , then with probability  $1/2$  sample  $u^{(1)}, \dots, u^{(N)} \stackrel{\text{iid}}{\sim} \ell_{\vartheta}(\cdot)$  and compute

$$r_u^N(\theta, \vartheta) = \frac{q(\vartheta, \theta)}{q(\theta, \vartheta)} \frac{\eta(\vartheta)}{\eta(\theta)} \frac{g_{\vartheta}(y)}{g_{\theta}(y)} \frac{1}{N} \sum_{i=1}^N \frac{g_{\theta}(u^{(i)})}{g_{\vartheta}(u^{(i)})},$$

or (i.e., with probability  $1/2$ ) sample  $u^{(1)} \sim \ell_{\vartheta}(\cdot)$  and  $u^{(2)}, \dots, u^{(N)} \stackrel{\text{iid}}{\sim} \ell_{\theta}(\cdot)$ , and compute  $r_u^N(\vartheta, \theta)$ . The interpretation of what MHAAR achieves in this particularly simple scenario

is transparent: when  $c = 1$  the right hand side average is a consistent estimator of  $C_\theta/C_\vartheta$ , suggesting that the algorithm can approximate the algorithm we would have liked to implement initially. The simplicity of this scenario, where the latent variable  $z$  is absent, also allows for a simple analysis illustrating the theoretical benefits of Algorithm 1. Establishing these results in full generality requires the use of convex order tools as in [Andrieu and Vihola, 2016], which is far beyond the scope of this paper. Instead performance improvement will be illustrated through numerical experiments.

Consider standard performance measures associated to a Markov transition probability  $\Pi$  of invariant distribution  $\nu$  defined on some measurable space  $(\mathbf{E}, \mathcal{E})$ . Let  $L^2(\mathbf{E}, \nu) := \{f: \mathbf{E} \rightarrow \mathbb{R}, \text{var}_\nu(f) < \infty\}$  and  $L_0^2(\mathbf{E}, \nu) := L^2(\mathbf{E}, \nu) \cap \{f: \mathbf{E} \rightarrow \mathbb{R}, \mathbb{E}_\nu(f) = 0\}$ . For any  $f \in L^2(\mathbf{E}, \nu)$  the asymptotic variance is defined as

$$\text{var}(f, \Pi) := \lim_{M \rightarrow \infty} \text{var}_\nu \left( M^{-1/2} \sum_{i=1}^M f(X_i) \right),$$

which is guaranteed to exist for reversible Markov chains (although it may be infinite) and for a  $\nu$ -reversible kernel  $\Pi$  its right spectral gap

$$\text{Gap}_R(\Pi) := \inf \{ \mathcal{E}_\Pi(f) : f \in L_0^2(\mathbf{E}, \nu), \text{var}_\nu(f) = 1 \},$$

where for any  $f \in L^2(\mathbf{E}, \nu)$   $\mathcal{E}_\Pi(f) := \frac{1}{2} \int_{\mathbf{E}} \nu(dx) \Pi(x, dy) [f(x) - f(y)]^2$  is the so-called Dirichlet form. The right spectral gap is particularly useful in the situation where  $\Pi$  is a positive operator, in which case  $\text{Gap}_R(\Pi)$  is related to the geometric rate of convergence of the Markov chain.

Hereafter we let  $\mathring{P}^N(\theta, d\vartheta)$  be the Markov chain transition kernel corresponding to Algorithm 1 in the absence of  $z$ .

**Theorem 2.** *With  $P$  and  $\mathring{P}^N$  as defined in (2) and corresponding to Algorithm 1, respectively,*

1. *for all  $N$ ,  $\text{Gap}_R(\mathring{P}^N) \leq \text{Gap}_R(P)$  and  $N \mapsto \text{Gap}_R(\mathring{P}^N)$  is non decreasing,*
2. *for any  $f \in L^2(\mathbf{X}, \pi)$ ,*
  - (a)  *$N \mapsto \text{var}(f, \mathring{P}^N)$  is non increasing,*
  - (b) *for all  $N$ ,  $\text{var}(f, \mathring{P}^N) \geq \text{var}(f, P)$ .*

The proof can be found in Appendix A.2. This result motivates the practical usefulness of the algorithm, in particular in a parallel computing environment. Indeed, one crucial property of Algorithm 1 is that for both updates  $Q_1^N(\cdot)$  and  $Q_2^N(\cdot)$ , sampling of  $u^{(1)}, \dots, u^{(N)}$  and computation of  $r_{u^{(1)}}(\theta, \vartheta), \dots, r_{u^{(N)}}(\theta, \vartheta)$  can be performed in a parallel fashion therefore opening the possibility to improve on the variance  $\text{var}(f, \mathring{P})$  of estimators, but more significantly the burn-in period of algorithms. Indeed one could object that running  $M \in \mathbb{N}^+$  independent chains in parallel with  $N = 1$  and combining their averages, instead of using the output from a single chain with  $N = M$  would achieve variance reduction. However our point is that the former does not speed up convergence to equilibrium, while the latter will, in general. Unfortunately, while estimating the asymptotic variance  $\text{var}(f, \mathring{P}^N)$  from simulations is achievable, estimating time to convergence to equilibrium is far from standard in general. The following toy example is an exception and illustrates our point.

**Example 4.** Here we let  $\pi$  be the uniform distribution on  $\Theta = \{-1, 1\}$ ,  $\mathbf{U} = \{a, a^{-1}\}$  for  $a > 0$ ,  $Q_{\theta, -\theta}(u = a) = 1/(1 + a)$ ,  $Q_{\theta, -\theta}(u = 1/a) = a/(1 + a)$  and

$$\varphi(\theta, \vartheta, \mathbf{u}, k, c) = (\vartheta, \theta, 1/u^{(1)}, u^{(2)}, \dots, u^{(N)}, k, 3 - c).$$

In other words  $\dot{P}$  can be reparametrised in terms of  $a$  and with the choice  $q(\theta, -\theta) = 1 - \alpha$  for  $\alpha \in [0, 1)$  we obtain

$$\dot{P}(\theta, -\theta) = (1 - \alpha) \left[ \frac{1}{1 + a} \min\{1, a\} + \frac{a}{1 + a} \min\{1, a^{-1}\} \right].$$

Note that there is no need to be more specific than say  $Q_{\theta, \theta}(u) > 0$  for  $(\theta, u) \in \mathbf{X} \times \mathbf{U}$  as then a proposed “stay” is always accepted. Now for  $N \geq 2$  and  $\theta \in \Theta$  we have

$$\begin{aligned} \dot{P}^N(\theta, -\theta) = \frac{1 - \alpha}{2} & \left[ \sum_{k=0}^N \beta^N(k) \min\{1, w_k(N)\} \right. \\ & \left. + \sum_{k=0}^N \left( \frac{a}{1 + a} \beta^{N-1}(k - 1) + \frac{1}{1 + a} \beta^{N-1}(k) \right) \min\{1, w_k^{-1}(N)\} \right], \end{aligned}$$

where  $\beta^N(k)$  is the probability mass function of the binomial distribution of parameters  $N$  and  $1/(1 + a)$  and  $w_k(N) := ka/N + (1 - k/N)a^{-1}$ . The second largest eigenvalue of the corresponding Markov transition matrix is  $\lambda_2(N) = 1 - 2\dot{P}^N(\theta, -\theta)$  from which we find the relaxation time  $T_{\text{relax}}(N) := 1/(2\dot{P}^N(\theta, -\theta))$ , and bounds on the mixing time  $T_{\text{mix}}(\epsilon, N)$ , that is the number of iterations required for the Markov chain to have marginal distribution within  $\epsilon$  of  $\pi$ , in the total variation distance, [Levin and Peres \[2017, Theorem 12.3 and Theorem 12.4\]](#)

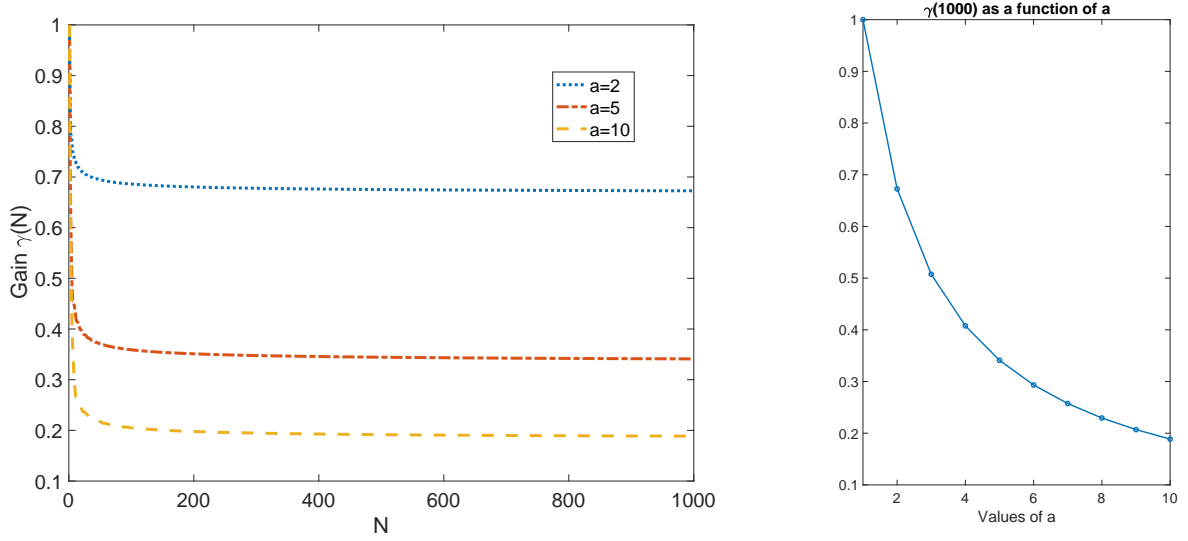
$$-(T_{\text{relax}}(N) - 1) \log(2\epsilon) \leq T_{\text{mix}}(\epsilon, N) \leq -T_{\text{relax}}(N) \log(\epsilon/2).$$

We define the relative burn-in time fraction,  $\gamma(N) := T_{\text{relax}}(N)/T_{\text{relax}}(1)$ , which is independent of  $\alpha$  and captures the benefit of MHAAR in terms of convergence to equilibrium. In [Figure 1](#) we present the evolution of  $N \mapsto \gamma(N)$  for  $a = 2, 5, 10$  and  $\gamma(1000)$  as a function of  $a$ . As expected the worse the algorithm corresponding to  $\dot{P}$  is, the more beneficial averaging is: for  $a = 2, 5, 10$  we observe running time reductions of approximately 35%, 65% and 80% respectively. This suggests that computationally cheap, but possibly highly variable, estimators of the acceptance ratio may be preferable to reduce burn-in when a parallel machine is available and communication costs are negligible.

## 2.5 Example: improving transdimensional samplers

The following example motivates the scenario considered in this section and on which we illustrate the interest of the proposed approach.

**Example 5** (Poisson multiple change-point model). The UK coal-mining disasters dataset consists of  $n$  records  $y_{1:n}$  of the number of disasters at a given set of dates. In [\[Green, 1995\]](#), it is proposed to model the dataset with a non-homogenous Poisson process model on the time interval  $[0, L]$  with a step-wise constant intensity function with changepoints



**Figure 1:** Evolution of  $N \mapsto \gamma(N)$  for  $a = 2, 5, 10$  and  $\gamma(1000)$  as a function of  $a$ .

$0 = s_0 < s_1 \dots < s_m = L$  and heights  $h_1, \dots, h_m \in \mathbb{R}_+$  for some  $m \in \mathbb{N}$ . Letting  $z_m := (\{s_j\}_{j=0}^m, \{h_j\}_{j=1}^m)$  the data likelihood under ‘model’  $m$  is therefore

$$\log \mathcal{L}_m(y_{1:n}; z_m) = \sum_{j=1}^m \log h_j \left( \sum_{i=1}^n \mathbb{I}_{[s_{j-1}, s_j)}(y_i) \right) - \sum_{j=1}^m h_j (s_j - s_{j-1}),$$

and inferring  $(m, z_m)$  is of interest. In a Bayesian framework one can ascribe a prior to  $(m, z_m)$  and infer both model and within model parameters from the associated posterior distribution. Sampling from such transdimensional distribution requires the use of a particular type of MH update, as proposed in [Green \[1995\]](#). Such algorithms may be difficult to design and we show how they can benefit from our approach.

In this section we consider target distributions  $\pi(\theta, dz_\theta)$  on  $\cup_{\vartheta \in \Theta} \{\vartheta\} \times Z_\vartheta$ , where in general  $\Theta \subseteq \mathbb{N}$  and the dimension  $d_\theta$  of  $Z_\theta \subset \mathbb{R}^{d_\theta}$  depends on  $\theta$ . We assume that  $\pi(\theta, dz_\theta)$  admits a density  $\pi(\theta, z_\theta)$  known up to a normalising constant, where  $z_\theta$  is a within model parameter. When sampling from this distribution a particular challenge is to define transdimensional transitions from  $(\theta, z_\theta)$  to  $(\vartheta, z_\vartheta)$  in situations where  $d_\theta \neq d_\vartheta$  and we focus on such updates only here. Practical algorithms consists of mixtures of such updates and more traditional within model updates [Green \[1995\]](#). Our aim here is to outline the solution proposed by [Green \[1995\]](#) and show how it fits, up to minor modifications, in the framework outlined in Subsection 2.3 and can benefit from the MHAAR methodology.

The main idea of [Green \[1995\]](#) consists, for  $\theta, \vartheta \in \Theta$ , of augmenting the within model parameters to ensure dimension matching, that is  $(z_\theta, u_{\theta, \vartheta}) \in Z_{\theta, \vartheta} := Z_\theta \times U_{\theta, \vartheta}$ ,  $(z_\vartheta, u_{\vartheta, \theta}) \in Z_{\vartheta, \theta} := Z_\vartheta \times U_{\vartheta, \theta}$ , with  $U_{\theta, \vartheta} \subset \mathbb{R}^{d_{\theta, \vartheta}}$ ,  $U_{\vartheta, \theta} \subset \mathbb{R}^{d_{\vartheta, \theta}}$  for  $d_{\theta, \vartheta}, d_{\vartheta, \theta} \in \mathbb{N}$  such that  $d_\theta + d_{\theta, \vartheta} = d_\vartheta + d_{\vartheta, \theta}$ , and defining extended distributions

$$\pi(\theta, d(z_\theta, u_{\theta, \vartheta})) = \pi(\theta, dz_\theta) Q_{\theta, \vartheta, z_\theta}(du_{\theta, \vartheta}),$$

for some probability distribution  $Q_{\theta, \vartheta, z_\theta}(\cdot)$  on  $U_{\theta, \vartheta}$ . Note that in some scenarios we may have  $d_{\theta, \vartheta} = 0$  (resp. or  $d_{\vartheta, \theta} = 0$ ), in which case  $u_{\theta, \vartheta}$  (resp.  $u_{\vartheta, \theta}$ ) should be ignored.



**Example 6** (Poisson multiple change-point model (ctd.)). For the coal-mining disaster a transdimensional update may consist of adding or removing a changepoint and its height, in which case  $u_{m,m+1} = (s_*, h_*, j) \in [0, L] \times \mathbb{R}_+ \times \llbracket m+1 \rrbracket$  and  $u_{m,m-1} \in \llbracket m \rrbracket$ . A possible choice for the distributions is the uniform distribution for  $u_{m,m-1}$  (a randomly chosen changepoint is removed) the uniform distribution for  $s_*$ , and the prior distribution for  $h_*$ , in which case  $j$  is a deterministic function of  $s_*$  and  $z_m$ . This update is referred to as ‘birth-death’.

Together with an invertible mapping  $\phi_{\theta,\vartheta} : Z_{\theta,\vartheta} \rightarrow Z_{\vartheta,\theta}$  such that  $\phi_{\theta,\vartheta}^{-1} = \phi_{\vartheta,\theta}$  this allows one to define the involution  $\varphi(\theta, \vartheta, z_\theta, u_{\theta,\vartheta}) = (\vartheta, \theta, z_\vartheta, u_{\vartheta,\theta})$  with  $(z_\vartheta, u_{\vartheta,\theta}) = \phi_{\theta,\vartheta}(z_\theta, u_{\theta,\vartheta})$ , and hence a valid MH type update [Green \[1995\]](#). While the choice of  $U_{\theta,\vartheta}$  and  $\phi_{\theta,\vartheta}$  are often natural for numerous problems, choosing the distribution  $Q_{\theta,\vartheta}$  can be difficult and result in poor performance. Our aim here is to show that averaging acceptance ratios of the standard procedure over multiple matching variables can improve performance significantly. The MHAAR algorithm in this context, which we call Reversible-multiple-jump MCMC (RmJ-MCMC), follows along the lines of Subsection 2.3.

The RmJ-MCMC update is described in detail in Algorithm 2 where we have taken into account Remark 2 and changed the distribution of  $c$ , but also taken into account that for  $c = 2$  the nature of the auxiliary variables may differ for  $i = k$  and  $i \neq k$ . In Algorithm 2,  $r_u(\theta, \vartheta, z_\theta)$  is the acceptance rate of the standard RJ-MCMC algorithm when the current sample is  $(\theta, z_\theta)$ ,  $\vartheta$  is proposed from  $q(\theta, \cdot)$  and  $u_{\theta,\vartheta}$  is the dimension-matching variable sampled from  $Q_{\theta,\vartheta,z_\theta}(\cdot)$ . One can check that Algorithm is a special case of MHAAR given in Algorithm 1, where the space of the latent variable depends on  $\theta$ , and likewise the space of auxiliary variables, which are the dimension-matching variables of RmJ-MCMC, depends on  $\theta, \vartheta$  as well as  $c$ .

---

**Algorithm 2:** RmJ-MCMC: MHAAR for trans-dimensional models.

---

**Input:** Current sample  $(\theta, z_\theta)$   
**Output:** New sample

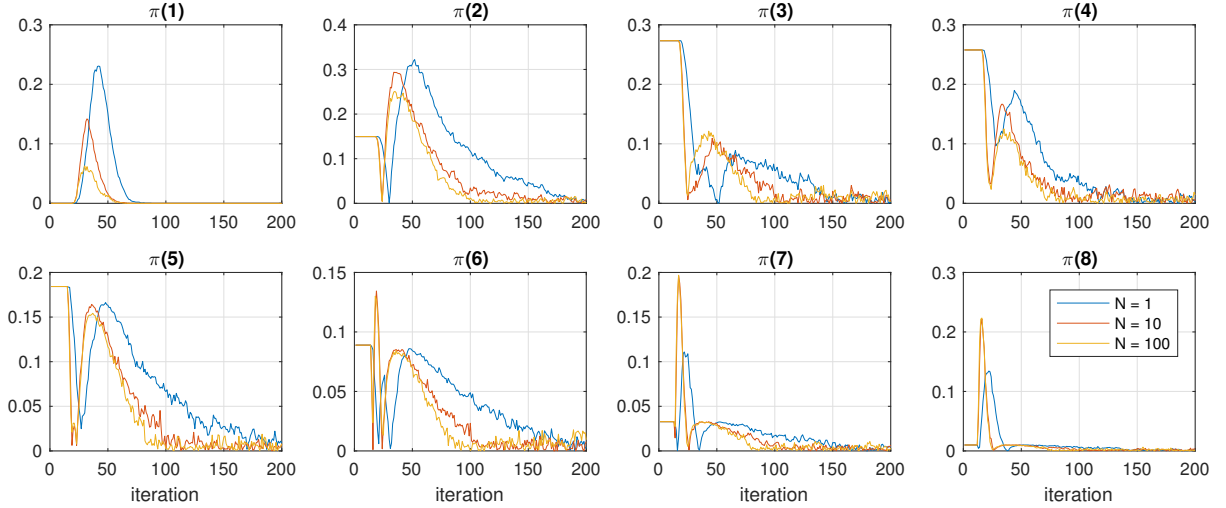
- 1 Sample  $\vartheta \sim q(\theta, \cdot)$  and  $c \sim \mathcal{P}(\omega(\theta, \vartheta, z_\theta, 1), \omega(\theta, \vartheta, z_\theta, 2))$ .
- 2 **if**  $c = 1$  **then**
- 3     **for**  $i = 1, \dots, N$  **do**
- 4         Sample  $u_{\theta,\vartheta}^{(i)} \sim Q_{\theta,\vartheta,z_\theta}(\cdot)$
- 5         Sample  $k \sim \mathcal{P}(r_{u_{\theta,\vartheta}^{(1)}}(\theta, \vartheta, z_\theta), \dots, r_{u_{\theta,\vartheta}^{(N)}}(\theta, \vartheta, z_\theta))$ , set  $z'_\vartheta = \phi_{\theta,\vartheta}^{[1]}(z_\theta, u_{\theta,\vartheta}^{(k)})$ .
- 6         Return  $(\vartheta, z'_\vartheta)$  with probability  $\min\{1, r_u^N(\theta, \vartheta, z_\theta)\}$ , otherwise return  $(\theta, z_\theta)$ .
- 7 **else**
- 8     Sample  $k \sim \text{Unif}(\llbracket N \rrbracket)$  and  $u_{\theta,\vartheta}^{(k)} \sim Q_{\theta,\vartheta,z_\theta}(\cdot)$ , set  $z'_\vartheta = \phi_{\theta,\vartheta}^{[1]}(z_\theta, u_{\theta,\vartheta}^{(k)})$ .
- 9     **for**  $i = 1, \dots, N, i \neq k$  **do**
- 10         Sample  $u_{\vartheta,\theta}^{(i)} \sim Q_{\vartheta,\theta,z_\vartheta}(\cdot)$ .
- 11     Return  $(\vartheta, z'_\vartheta)$  with probability  $\min\{1, r_u^N(\vartheta, \theta, z'_\vartheta)^{-1}\}$ , otherwise return  $(\theta, z_\theta)$ .

---

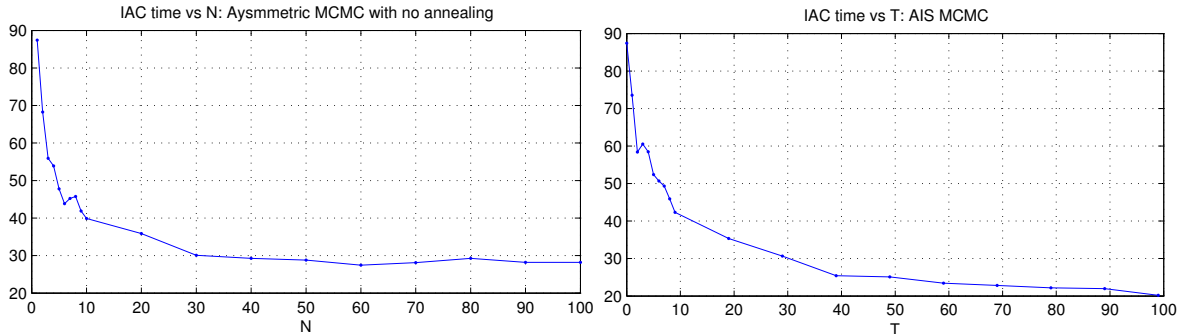
**Example 7** (Poisson multiple change-point model (ctd.)). We now evaluate this approach on the coal-mining disaster example. In order to improve computational efficiency



we set  $\omega(m, m+1, 1) = 1$  and  $\omega(m, m-1, 1) = 0$ . Indeed when attempting a birth it is preferable to average over the continuous valued  $u_{m,m+1}$  rather than the discrete valued  $u_{m,m-1}$ , in particular when  $N \gg m+1$ . The priors chosen are as in [Green \[1995\]](#) and the specifics of the MCMC move for updating the latent variables within model are chosen as in [Karagiannis and Andrieu \[2013\]](#). To illustrate the gains in terms of convergence to equilibrium of our scheme we had 3000 independent runs started at the same point  $x_0$ , estimated the expectations  $\mathbb{E}^N[\mathbb{I}\{M_t = m\}]$  by an ensemble average, and reported  $|\hat{\pi}(m) - 3000^{-1} \sum_{k=1}^{3000} \mathbb{I}\{M_t^{(k)} = m\}|$  for  $m \in \{1, \dots, 8\}$  and  $N = 1, 10, 100$  in [Figure 2](#) where  $\hat{\pi}(m)$  was estimated by a realisation of length  $10^6$  with  $N = 90$  and  $T = 50$ , discarding the burn-in. We see that the approach appears to reduce time to convergence to equilibrium by the order of 50%. We also generated  $K = 10^6$  samples to compute the IAC for  $m$ . [Figure 3](#) indicates a variance reduction of the order of 60% at  $N = 130$ . We also provide results for the scheme used in [Karagiannis and Andrieu \[2013\]](#) (referred to as AIS for Annealing Importance Sampling) for illustration. For  $T$  (a tuning parameter of the algorithm) large the algorithm approaches the algorithm which would sample from the model distribution directly as  $T \rightarrow \infty$ . Our algorithm achieves similar performance improvement, but is parallelisable.



**Figure 2:** Estimates of time to convergence of  $\mathbb{E}_{x_0}^N[f_m(X_i)]$  to  $\pi(m)$  for  $N = 1, 10, 100$ .



**Figure 3:** Left: IAC for  $m$  vs number of particles  $N = 1, 2, \dots, 10, 20, \dots, 100$ . Right: IAC for  $m$  vs number of particles  $T = 0, 1, 2, \dots, 10, 20, \dots, 100$ .

### 3 An efficient alternative to pseudo-marginal algorithms

In this section we consider a class of latent variable models of probability density defined on  $\Theta \times \mathcal{Z}^T$  for some  $T \geq 1$  and of the form

$$\pi(\theta, z) \propto \eta(\theta) \prod_{t=1}^T \gamma_{t,\theta}(z_t),$$

with  $z = z_{1:T}$ ,  $\eta(\theta)$  a prior density, and  $\gamma_{t,\theta}(z_t)$  is typically a complete likelihood function depending on some observation  $y_t$ , see the examples in this section. We drop any such dependencies from notation for simplicity. It was shown in [Yildirim et al. \[2018\]](#) that it is possible to develop efficient sampling schemes for such models which in particular scale favourably with  $T$  large. We show here that these algorithms can be further improved at little cost using the methodology developed in this paper, leading in particular to alternative to pseudo-marginal algorithms [Andrieu and Roberts \[2009\]](#) with much better performance. We will show in Section 4 how these ideas can be extended to the context of state-space models.

#### 3.1 A novel consistent pseudo-marginal estimator

The algorithm we develop can be thought of as being the numerical approximation of the scheme in Subsection 2.2 where  $u = z'_{1:T}$  the proposed new values of the MH update and  $Q_{\theta,\vartheta,z}$  has density proportional to  $\prod_{t=1}^T \gamma_{t,\theta,\vartheta}(z'_t)$ . This cannot be achieved in practice and instead replace this update with a Markov kernel reversible with respect to this distribution, in the spirit of [Neal \[2004\]](#), dependent on a parameter  $M \in \mathbb{N}_*$  and such that as  $M \rightarrow \infty$  the algorithm approaches the idealised algorithm. An interesting feature is that this kernel produces multiple samples which can be used in the averaging procedure, at very little extra cost.

We first introduce the algorithm of [Yildirim et al. \[2018\]](#) and identify computational inefficiencies which can be addressed with a MHAAR strategy. For  $\theta, \vartheta \in \Theta$  and  $t \in \llbracket T \rrbracket$ , let  $q_{t,\theta,\vartheta}$  be a probability distribution on  $(\mathcal{Z}, \mathcal{Z})$  and with  $M \geq 2$  and  $\mathbf{u} := u_{1:T}^{(1:M-1)} \in \mathcal{Z}^{(M-1)T}$ , let

$$\Phi_{\theta,\vartheta}(\mathrm{d}\mathbf{u}) := \prod_{t=1}^T \prod_{i=1}^{M-1} q_{t,\theta,\vartheta}(\mathrm{d}u_t^{(i)}). \quad (17)$$

For notational simplicity we define  $\mathbf{v} = v_{1:T}^{(1:M)} \in \mathcal{Z}^{MT}$  such that  $(v_1^{(1)}, \dots, v_T^{(1)}) = z$  and  $v_{1:T}^{(2:M)} = \mathbf{u}$ . For any index sequence  $\mathbf{k} = (k_1, \dots, k_T) \in \llbracket M \rrbracket^T$ , we define  $v^{(\mathbf{k})} := (v_1^{(k_1)}, \dots, v_T^{(k_T)})$ , so that  $z = v^{(\mathbf{1})}$  where  $\mathbf{1} := (1, \dots, 1)$  is the vector of size  $T$  consisting of 1's. The proposal mechanism of the algorithms considered consists of sampling candidates from  $\Phi_{\theta,\vartheta}(\mathrm{d}\mathbf{u})$  and then attempting a swap of  $v^{(\mathbf{k})}$  and  $v^{(\mathbf{j})}$ , where  $\mathbf{k} \in \llbracket M \rrbracket^T$  is sampled conditional on  $\mathbf{v}$  according to

$$b_{\theta,\vartheta}(\mathbf{k}|\mathbf{v}) := \prod_{t=1}^T \frac{\gamma_{t,\theta,\vartheta}(v_t^{(k_t)})/q_{t,\theta,\vartheta}(v_t^{(k_t)})}{\sum_{j=1}^M \gamma_{t,\theta,\vartheta}(v_t^{(j)})/q_{t,\theta,\vartheta}(v_t^{(j)})}. \quad (18)$$

Here  $\gamma_{t,\theta,\vartheta}(v)$  is a user defined probability density on  $(Z, \mathcal{Z})$ —possible choices include  $\gamma_{t,(\theta+\vartheta)/2}(v)$  or  $\gamma_{t,\theta}(v)$ . Using the framework of Subsection 2.1 we let  $\xi = (\theta, \vartheta, \mathbf{v}, \mathbf{k})$ ,

$$\pi(d\xi) := \pi(d(\theta, z))q(\theta, d\vartheta)\Phi_{\theta,\vartheta}(d\mathbf{u})b_{\theta,\vartheta}(\mathbf{k}|\mathbf{v}),$$

and consider the involution  $\varphi(\theta, \vartheta, \mathbf{v}, \mathbf{k}) = (\vartheta, \theta, \mathbf{s}_{1,\mathbf{k}}(\mathbf{v}), \mathbf{k})$  where  $\mathbf{s}_{1,\mathbf{k}} : Z^{MT} \mapsto Z^{MT}$  is the operator on  $\mathbf{v}$  which swaps  $v^{(1)}$  and  $v^{(\mathbf{k})}$ , that is, if  $\mathbf{v}' = \mathbf{s}_{1,\mathbf{k}}(\mathbf{v})$ , it satisfies

$$\mathbf{v}'_t = \begin{cases} v_t^{(1)} & \text{for } i = k_t, \\ v_t^{(k_t)} & \text{for } i = 1, \dots, n; i = 1, \dots, M. \\ v_t^{(i)} & \text{otherwise.} \end{cases} \quad (19)$$

The corresponding acceptance ratio can be shown to be  $r_{v^{(1)},v^{(\mathbf{k})}}(\theta, \vartheta)$ , where, for any  $z, z' \in Z^T$ ,

$$r_{z,z'}(\theta, \vartheta) := \frac{q(\vartheta, \theta)\eta(\vartheta)}{q(\theta, \vartheta)\eta(\theta)} \prod_{t=1}^T \frac{\gamma_{t,\theta,\vartheta}(z_t)}{\gamma_{t,\theta}(z_t)} \frac{\gamma_{t,\vartheta}(z'_t)}{\gamma_{t,\theta,\vartheta}(z'_t)}, \quad (20)$$

We will refer to this algorithm as AIS MCMC, since the proposal mechanism for  $z$  can be viewed as a one-step annealing using the ‘intermediate’ distribution with (unnormalised) density  $\gamma_{t,\theta,\vartheta}(\cdot)$ , building on the ideas in Neal [2004]. While this algorithm can be shown to be efficient in the regime  $T \rightarrow \infty$  by appropriate scaling of  $\vartheta - \theta$ , it should be clear that the use of a single “path”  $\mathbf{k}$  in  $\mathbf{v}$  is wasteful and the use of the “Rao-Blackwellised” acceptance ratio

$$r_{1,\mathbf{v}}(\theta, \vartheta) := \sum_{\mathbf{k} \in \llbracket M \rrbracket^T} b_{\theta,\vartheta}(\mathbf{k}|\mathbf{v}) r_{v^{(1)},v^{(\mathbf{k})}}(\theta, \vartheta), \quad (21)$$

may be preferable. Before showing how this can be achieved within the MHAAR framework we take a closer look at  $r_{1,\mathbf{v}}(\theta, \vartheta)$ , which further motivates these algorithms. Rearranging terms (see Theorem 2.3) it can be shown that

$$r_{1,\mathbf{v}}(\theta, \vartheta) = \frac{q(\vartheta, \theta)\eta(\vartheta)}{q(\theta, \vartheta)\eta(\theta)} \prod_{t=1}^T \frac{\gamma_{t,\theta,\vartheta}(v_t^{(1)})}{\gamma_{t,\theta}(v_t^{(1)})} \prod_{t=1}^T \frac{\sum_{i=1}^M \gamma_{t,\vartheta}(v_t^{(i)})/q_{t,\theta,\vartheta}(v_t^{(i)})}{\sum_{j=1}^M \gamma_{t,\theta,\vartheta}(v_t^{(j)})/q_{t,\theta,\vartheta}(v_t^{(j)})}, \quad (22)$$

implying in particular that this can be computed in  $\mathcal{O}(MT)$  operations and not  $\mathcal{O}(M^T)$  as suggested by our earlier expression. It is worth noting that for any  $\theta, \vartheta \in \Theta$ , this is an unbiased estimator of  $r(\theta, \vartheta)$  when  $z = v^{(1)} \sim \pi(dz | \theta)$  – this is established in a more general context in Theorem 4 in Section 4.2. The choice  $\gamma_{t,\theta,\vartheta} = \gamma_{t,\theta}$  leads to

$$r_{1,\mathbf{v}}(\theta, \vartheta) = \frac{q(\vartheta, \theta)\eta(\vartheta)}{q(\theta, \vartheta)\eta(\theta)} \prod_{t=1}^T \frac{\sum_{i=1}^M \gamma_{t,\vartheta}(v_t^{(i)})/q_{t,\theta,\vartheta}(v_t^{(i)})}{\sum_{j=1}^M \gamma_{t,\theta}(v_t^{(j)})/q_{t,\theta,\vartheta}(v_t^{(j)})}, \quad (23)$$

which is reminiscent of the acceptance ratio of a pseudo-marginal algorithm Andrieu and Roberts [2009] where importance sampling is used to estimate the likelihood function. However the crucial difference here is that only one set of auxiliary variables, sampled afresh at each iteration, is used to estimate the numerator and denominator of  $r(\theta, \vartheta)$  in (1), leading to reduced variability and improved performance – as pointed out in Subsection 1.2, for a pseudo-marginal algorithm a poor draw of the denominator leads

to the algorithm getting stuck in the same state for a large number of iterations. This algorithm can be thought of as an alternative to the correlated pseudo-marginal algorithm of Deligiannidis et al. [2018].

The new algorithm, MHAAR-RB (for Rao-Blackwellised) hereafter, is obtained by alternating between two sampling mechanisms for  $\mathbf{k}$ . Let  $\xi = (\theta, \vartheta, \mathbf{v}, \mathbf{k}, c) \in \Theta^2 \times \mathbb{Z}^{MT} \times \llbracket M \rrbracket^T \times \{1, 2\}$  and

$$\hat{\pi}(\mathrm{d}\xi) := \frac{1}{2} \pi(\mathrm{d}(\theta, z)) Q_c^M((\theta, z); \mathrm{d}(\mathbf{u}, \mathbf{k})) \quad (24)$$

with

$$\begin{aligned} Q_1^M((\theta, z); \mathrm{d}(\mathbf{u}, \mathbf{k})) &:= q(\theta, \mathrm{d}\vartheta) \Phi_{\theta, \vartheta}(z, \mathrm{d}\mathbf{u}) b_{\theta, \vartheta}^{(1)}(\mathbf{k}|\mathbf{v}), \\ Q_2^M((\theta, z); \mathrm{d}(\mathbf{u}, \mathbf{k})) &:= q(\theta, \mathrm{d}\vartheta) \Phi_{\vartheta, \theta}(z, \mathrm{d}\mathbf{u}) b_{\vartheta, \theta}^{(2)}(\mathbf{k}|\mathbf{v}). \end{aligned}$$

where the sampling probabilities are given as

$$b_{\theta, \vartheta}^{(1)}(\mathbf{k}|\mathbf{v}) := \prod_{t=1}^T \frac{\gamma_{t, \vartheta}(v^{(k_t)})/q_{t, \theta, \vartheta}(v^{(k_t)})}{\sum_{j=1}^M \gamma_{t, \vartheta}(v^{(j)})/q_{t, \theta, \vartheta}(v^{(j)})}, \quad (25)$$

which can be shown to be obtained by weighting  $b_{\theta, \vartheta}(\mathbf{k}|\mathbf{v})$  by the acceptance ratio  $\hat{r}_{\mathbf{l}, \mathbf{k}, \mathbf{v}}(\theta, \vartheta)$  corresponding to  $\mathbf{k}$ , and  $b_{\theta, \vartheta}^{(2)}(\mathbf{k}|\mathbf{v}) = b_{\vartheta, \theta}(\mathbf{k}|\mathbf{v})$ . Given the current sample  $(\theta, z) \in \Theta \times \mathbb{Z}$ , an update of MHAAR-RB proceeds as follows:

1. Sample  $c \sim \text{Unif}(\{1, 2\})$ ,  $(\vartheta, \mathbf{u}, \mathbf{k}) \sim Q_c^N((\theta, z); \cdot)$  and form  $\xi = (\theta, \vartheta, \mathbf{v}, \mathbf{k}, c)$ .
2. With  $\mathfrak{s}_{\mathbf{l}, \mathbf{k}}(\mathbf{v})$  as in (19), let

$$\xi' = \varphi(\theta, \vartheta, \mathbf{v}, \mathbf{k}, c) := (\vartheta, \theta, \mathfrak{s}_{\mathbf{l}, \mathbf{k}}(\mathbf{v}), \mathbf{k}, 3 - c). \quad (26)$$

3. Return  $\xi'$  with probability  $\min\{1, \hat{r}(\xi)\}$ , otherwise return  $\xi$ , where

$$\hat{r}(\xi) := \begin{cases} r_{\mathbf{l}, \mathbf{v}}(\theta, \vartheta), & c = 1 \\ 1/r_{\mathbf{k}, \mathbf{v}}(\vartheta, \theta), & c = 2 \end{cases}, \quad (27)$$

with, for  $\mathbf{l} \in \llbracket M \rrbracket^T$ ,

$$r_{\mathbf{l}, \mathbf{v}}(\theta, \vartheta) := \frac{q(\vartheta, \theta) \eta(\vartheta)}{q(\theta, \vartheta) \eta(\theta)} \prod_{t=1}^T \frac{\gamma_{t, \theta, \vartheta}(v_t^{(l_t)})}{\gamma_{t, \theta}(v_t^{(l_t)})} \prod_{t=1}^T \frac{\sum_{i=1}^M \gamma_{t, \vartheta}(v_t^{(i)})/q_{t, \theta, \vartheta}(v_t^{(i)})}{\sum_{j=1}^M \gamma_{t, \theta, \vartheta}(v_t^{(j)})/q_{t, \theta, \vartheta}(v_t^{(j)})}. \quad (28)$$

The following theorem, whose proof is left to Appendix B.1, establishes the correctness of the acceptance ratio above.

**Theorem 3.** *The acceptance ratio resulting from the choices of  $\hat{\pi}$  as in (24) and the involution as in (26) is given by (27)-(28).*

A detailed pseudo-code of MHAAR-RB is given in Algorithm 3. When  $c = 1$ , the acceptance ratio does not depend on  $\mathbf{k}$ , which can be taken advantage of by sampling  $\mathbf{k}$

upon acceptance only. Notice also the optional stage which has not been discussed yet. These are motivated by the fact that the proposed variables  $(\vartheta, v^{(\mathbf{k})})$  are either accepted or rejected jointly and it seems natural, upon rejection, to attempt to refresh the current latent variable only, i.e. attempt a transition to  $(\theta, v^{(\mathbf{l})})$  for some  $\mathbf{l} \in \llbracket M \rrbracket^T$ . We show in Appendix B.2 that such a delayed rejection strategy is possible in general and takes the particular form shown in Algorithm 3, that is no rejection occurs in this optional stage in the situation where  $\gamma_{t,\theta,\vartheta} = \gamma_{t,\theta}$ . The computational cost of these steps is  $\mathcal{O}(MT)$ .

---

**Algorithm 3:** MHAAR-RB for the multiple latent variable model

---

**Input:** Current sample  $(\theta, z)$   
**Output:** New sample

- 1 Sample  $\vartheta \sim q(\theta, \cdot)$  and  $c \sim \text{Unif}(\{1, 2\})$
- 2 **if**  $c = 1$  **then**
  - 3 Set  $v^{(1)} = z_{1:T}$  and sample  $v_t^{(i)} \sim q_{t,\vartheta,\theta}(\cdot)$  for  $i = 2, \dots, M$ ,  $t = 1, \dots, T$ .
  - 4 Sample  $\mathbf{k} \sim b_{\theta,\vartheta}^{(1)}(\cdot|\mathbf{v})$  and set  $z' = v^{(\mathbf{k})}$ .
  - 5 Return  $(\vartheta, z')$  with probability  $\min\{1, r_{\mathbf{l},\mathbf{v}}(\theta, \vartheta)\}$ ; otherwise return  $(\vartheta, z)$ .
- 6 **else**
  - 7 Set  $v^{(1)} = z_{1:T}$  and sample  $v_t^{(i)} \sim q_{t,\vartheta,\theta}(\cdot)$  for  $i = 2, \dots, M$ ,  $t = 1, \dots, T$ .
  - 8 Sample  $\mathbf{k} \sim b_{\theta,\vartheta}^{(2)}(\cdot|\mathbf{v})$  and set  $z' = v^{(\mathbf{k})}$ .
  - 9 Return  $(\vartheta, z')$  with probability  $\min\{1, r_{\mathbf{k},\mathbf{v}}(\vartheta, \theta)^{-1}\}$ ; otherwise return  $(\theta, z)$ .
- 10 **Optional refreshment of**  $z_{1:T}$
- 11 **if** *the move is rejected and  $\gamma_{t,\theta,\vartheta} = \gamma_{t,\theta}$  for all  $t = 1, \dots, T$* , **then**
- 12 Sample  $\mathbf{l} \sim b_{\theta,\vartheta}^{\text{ref},(c)}(\cdot|\mathbf{v})$ , and set  $z_{1:T} = v^{(\mathbf{l})}$ , where
$$b_{\theta,\vartheta}^{\text{ref},(1)}(\mathbf{l}|\mathbf{v}) = \prod_{t=1}^n \frac{\gamma_{t,\theta}(v_t^{(\mathbf{l}_t)})/q_{t,\theta,\vartheta}(v_t^{(\mathbf{l}_t)})}{\sum_{i=1}^M \gamma_{t,\theta}(v_t^{(i)})/q_{t,\theta,\vartheta}(v_t^{(i)})}, \quad b_{\theta,\vartheta}^{\text{ref},(2)}(\mathbf{l}|\mathbf{v}) = \prod_{t=1}^n \frac{\gamma_{t,\theta}(v_t^{(\mathbf{l}_t)})/q_{t,\theta,\vartheta}(v_t^{(\mathbf{l}_t)})}{\sum_{i=1}^M \gamma_{t,\theta}(v_t^{(i)})/q_{t,\theta,\vartheta}(v_t^{(i)})}.$$
- 13 Return  $(\theta, z_{1:T})$ .

---

## 3.2 Examples

**Example 8** (ABC learning of an  $\alpha$ -stable distribution). Consider an intractable likelihood function  $\theta \mapsto \ell_\theta(y)$  for  $y \in \mathbf{Y}$  such that sampling from the corresponding data generating distribution is tractable. Approximate Bayesian computation (ABC) [Pritchard et al., 1999, Beaumont et al., 2002, Marjoram et al., 2003] is a general methodology to address inference in such scenarios. For  $\epsilon > 0$ , let  $g^\epsilon(z, y) := \kappa(y, z; \epsilon)$  for  $y, z \in \mathbf{Y}$ , where  $\kappa(\cdot; \cdot; \epsilon)$  is some kernel and  $\epsilon > 0$  a bandwidth parameter. An ABC-based approximation to the intractable posterior  $\pi(\theta) \propto \eta(\theta) \prod_{t=1}^T \ell_\theta(y_t)$  is obtained by marginalisation of the joint density

$$\pi^\epsilon(\theta, z_{1:T}) := \eta(\theta) \prod_{t=1}^T \ell_\theta(z_t) g^\epsilon(z_t, y_t). \quad (29)$$

For illustration we consider the scenario where the observations are assumed to arise from an  $\alpha$ -stable distribution  $\mathcal{A}(\alpha, \beta, \mu, \sigma)$ , where  $\alpha, \beta, \mu, \sigma \in \mathbb{R}_+$  are the shape, skewness, location, and scale parameters, respectively. Here we take  $g^\epsilon(z, y) = \kappa(\arctan(y); \arctan(z), \epsilon)$ , where  $\kappa(\cdot, \cdot; \epsilon)$  is taken a Gaussian kernel, as in [Yildirim et al. \[2015\]](#).

We generated a sequence of i.i.d. observations of length  $T = 100$  from  $\mathcal{A}(1.8, 0, 0, 2)$ . Assuming  $\beta$  is known, we consider estimating  $\theta = (\alpha, \mu, \sigma)$  using the ABC posterior distribution with  $\epsilon = 0.1$ . In order to illustrate the benefit of MHAAR-RB we compare performance of the RB and non-RB versions of the algorithms for two choices of  $\gamma_{t,\theta,\vartheta}(z)$ :

- $\gamma_{t,\theta,\vartheta}(z) = \ell_\theta(z)g^\epsilon(z, y_t)$  and  $Q_{t,\theta,\vartheta}(z) = \ell_\theta(z)$ . We refer to this version as MHAAR-RB-0 and the corresponding non-RB version is referred to as MwG (since the algorithm then corresponds to alternating between an update of  $z$  conditional upon  $\theta$  and  $\theta$  conditional upon  $z$ ),
- $\gamma_{t,\theta,\vartheta}(z) = \ell_{(\theta+\vartheta)/2}(z)g^\epsilon(z, y_t)$  and  $Q_{t,\theta,\vartheta}(z) = \ell_{(\theta+\vartheta)/2}(z)$ . We will refer to this version as MHAAR-RB-1. When no Rao-Blackwellisation is performed we refer to the algorithm as AIS MCMC.

Note that for a given  $M$  the complexity of these algorithms is comparable. The computational overhead arising from RB is limited since it consists of applying simple operations such as additions and multiplications to the most expensive quantities computed by all the algorithms. We provide precise details concerning prior choices and proposal distributions below and focus first on results.

We ran the algorithms for  $2 \times 10^5$  iterations for the values  $M = 10, 20, 50, 100$ . In [Table 4](#), we report IAC and  $\text{IAC} \times \text{CPU time per iteration}$  for the MHAAR-RB algorithms as well as their non-RB counterparts. The difference between the RB algorithms and their non-RB counterparts is striking: the former seem to benefit highly from increasing  $M$  in contrast with the latter. MHAAR-RB-0 seems superior to MHAAR-RB-1, which is explained by the fact the acceptance ratio of MHAAR-RB-0 in [\(23\)](#) enjoys a full averaging and suffers less from the dependency on  $z_t$  compared to the acceptance ratio of MHAAR-RB-1 in [\(22\)](#). We further observe the following further benefit of better mixing: for MHAAR-RB-0 the gain of using  $M = 100$  rather than  $M = 10$  replicas is  $1043/21 \approx 50$  while averaging the output from 10 computers running MHAAR-RB-0 for  $M = 10$  would have lead to a gain of 10. This advantage persists when (serial) CPU time is taken into account, even though our implementation uses Matlab, for which for loops can be particularly slow.

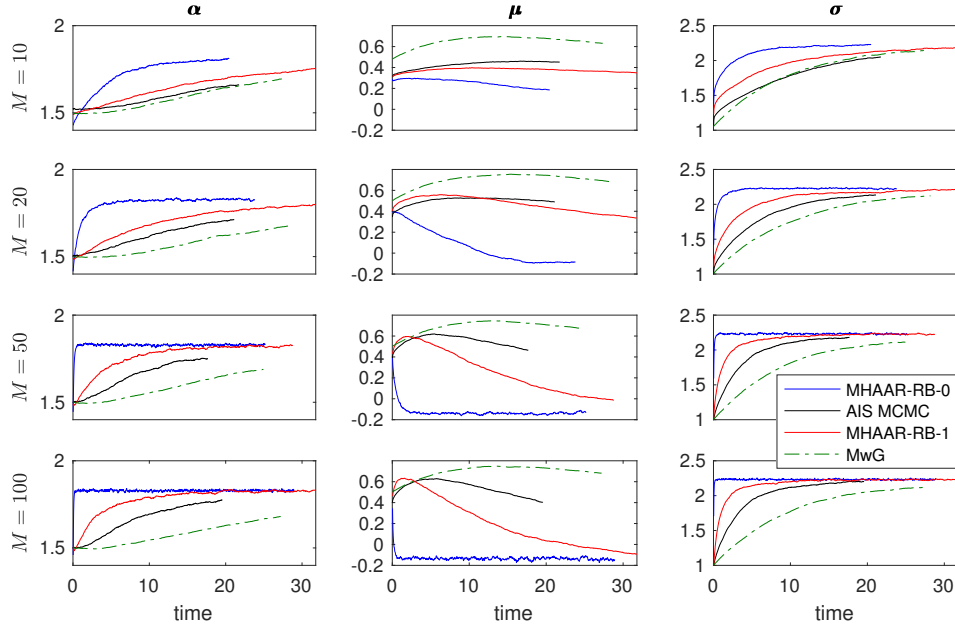
In [Figure 4](#) we report ensemble averages, over 1000 independent runs, vs time, for the algorithms compared in this example. One can observe the benefit of using averaging with MHAAR-RB, especially with the one without annealing, MHAAR-RB-0, as well as increasing  $M$ .

For all algorithms,  $\vartheta$  is proposed using a random walk proposal for all of its components, with standard deviation 0.2 for each. For simplicity, we take a flat prior for  $\theta$ . The first quarter of the  $2 \times 10^5$  iterations are discarded as burn-in time from the calculations related to IAC time.

**Example 9** (Gaussian process regression model). The Gaussian process regression model is an example for a single latent variable model, i.e.,  $T = 1$ . We observe pairs  $(x_i, y_i)$  for

$\theta$	$M$	IAC time				IAC $\times$ CPU time per iteration			
		MHAAR-RB-0	MwG	MHAAR-RB-1	AIS MCMC	MHAAR-RB-0	MwG	MHAAR-RB-1	AIS MCMC
$\alpha$	10	1043	1679	1504	1550	0.227	0.473	0.570	0.372
	20	244	1582	883	1198	0.088	0.635	0.465	0.407
	50	46	1748	250	1127	0.039	1.340	0.241	0.686
	100	21	1103	251	757	0.030	1.469	0.395	0.795
$\mu$	10	15952	11399	7519	2620	3.467	3.211	2.850	5.602
	20	2909	7706	3575	9235	1.054	3.093	1.880	3.135
	50	440	8041	2093	8254	0.375	6.166	2.022	5.024
	100	115	17323	2236	2665	0.165	23.068	3.524	2.8
$\sigma$	10	876	1461	2485	23317	0.190	0.412	0.942	0.629
	20	265	1175	489	1355	0.096	0.472	0.257	0.46
	50	74	990	243	959	0.063	0.759	0.235	0.584
	100	40	1160	257	576	0.057	1.545	0.405	0.605

**Table 1:** Comparison of algorithms in terms of IAC and IAC  $\times$  CPU time per iteration



**Figure 4:** Ensemble averages vs time for for MHAAR-RB-0, AIS MCMC, MHAAR-RB-1, MwG.

$i = 1, \dots, n$ , where  $x_i$  is a vector of  $d$  covariates, and

$$y_i = f(x_i) + \varepsilon_i, \quad \varepsilon_i \sim \mathcal{N}(0, \sigma^2),$$

where  $f$  is an unknown function with a Gaussian process prior with zero mean and some covariance function  $C$ , yielding  $(f(x_1), \dots, f(x_n)) \sim \mathcal{N}(0, C_{x_{1:n}})$ . One commonly used covariance function has the form

$$C_{x_{1:n}}(i, j) = \tau^2 \Upsilon(i, j), \quad \Upsilon(i, j) = \left( v + \exp \left\{ - \sum_{k=1}^d [\theta_i(x_i(k) - x_j(k))]^2 \right\} + \varsigma \delta_{i,j} \right)$$



We assume  $v$  and  $\varsigma$  are fixed and known, and the unknown variables  $\theta$  and  $z = (\tau, \sigma^2)$  are *a priori* independent, having Gaussian prior distributions for their logarithms. The log-likelihood of  $y = y_{1:n}$  and  $x = x_{1:n}$  given  $\theta$  is

$$\ell(\theta, z; x, y) = -0.5 (\log |\Sigma| + y^T \Sigma^{-1} y) .$$

where  $\Sigma = \tau^2 \Upsilon + \sigma^2 I_n$ . Therefore we have a joint distribution  $\pi(\theta, z)$  over the unknown variables.

Following the terminology of Neal [2004, 2010], we call a variable a slow (resp. fast) variable if the update of the posterior density is hard (resp. easy) when the variable is changed with the other parameters fixed. If eigenvalue decomposition is used for  $\Sigma$ , then  $\theta$  may be viewed as the slow variable, and  $z = (\tau, \sigma^2)$  as fast variables: Suppose  $\Upsilon = E \Lambda E^T$  and  $\hat{y} = E^T y$ , where  $\Lambda$  has the eigenvalues  $\lambda_1, \dots, \lambda_n$  on its diagonal. Since  $\Upsilon e = \lambda e$  implies  $\Sigma e = (\tau^2 \lambda + \sigma^2) e$ , we can write

$$\ell(y; \theta) = -0.5 \sum_{i=1}^n \log(\tau^2 \lambda_i + \sigma^2) - 0.5 \sum_{i=1}^n \frac{\hat{y}_i^2}{\tau^2 \lambda_i + \sigma^2} .$$

Therefore, while  $\hat{y}_i$ 's have to be re-evaluated when  $\theta$  changes, changing  $z = (\tau, \sigma^2)$  does not require re-evaluation of  $\hat{y}_i$ 's, which is the most computationally demanding part of the likelihood evaluation. The distinction of slow-fast variables for this model gets clearer for larger  $n$ . The fact that  $\theta$  is the slow variable and  $z$  is the fast variable justifies the use of MHAAR-RB, whose performance increases with the number of auxiliary variables generated for the latent variable  $z$ . When  $\vartheta$  is proposed, one needs to perform a single eigenvalue decomposition, which is expensive, which is followed by sampling  $M$  auxiliary variables and calculating quantities depending on them, which is relatively cheap even for large values of  $M$ .

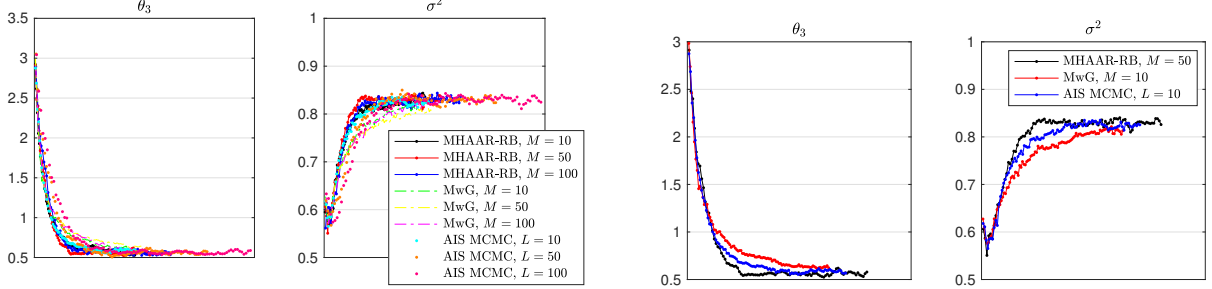
We have compared AIS MCMC with MHAAR-RB for the Gaussian process regression model on the same data set used in Neal [2010], with  $p = 12$  covariates and  $n = 100$  points. (The software in <https://www.cs.toronto.edu/~radford/ensmcmc.software.html> can be used to generate the data.) We ran MHAAR-RB with  $\gamma_{\theta, \vartheta} = \gamma_{\theta}$  with values  $M = 10, 50, 100$  and AIS MCMC with  $L = 10, 50, 100$  annealing steps with a geometric annealing schedule. All algorithms are started from the same initial point and run for  $10^6$  iterations. To demonstrate performance, we report IAC and  $\text{IAC} \times \text{CPU}$  times for the average  $\frac{1}{d} \sum_{i=1}^d \theta_i$  in Table 2. Although the difference between the performances of the two algorithms is not spectacular in terms of IAC, the MHAAR-RB algorithm benefits from its simplicity in generating the auxiliary variables, hence beating AIS MCMC significantly in terms of  $\text{IAC} \times \text{CPU}$  time. The table also shows the poorer performance of a MwG algorithm, where both slow and fast variables are updated in an alternating fashion by MH moves with random walk proposals. This indicates the usefulness of algorithms such as AIS MCMC and MHAAR-RB that exploit the existence of slow vs fast variables.

We also report some results pertaining the converge of the algorithms for this example. Figure 5 shows ensemble averages of the compared algorithms, out of 1000 runs starting from the same initial point, versus both iteration (left) and time (right). We ran the algorithms with all the parameter choices appearing in Table 2. The parameter choices appearing in the figure,  $M = 50$  for MHAAR-RB,  $M = 10$  for MwG and  $L = 10$  for AIS MCMC, correspond to the best choices in terms of convergence vs time. The figure



**Table 2:** IAC and IAC  $\times$  CPU times

$M$ or $L$	IAC time ( $\times 10^3$ )			IAC $\times$ CPU per iteration		
	MHAAR-RB	MwG	AIS MCMC	MHAAR-RB	MwG	AIS MCMC
10	2.33	5.12	2.24	15.98	32.32	15.89
50	1.06	4.21	1.37	8.54	30.6965	16.09
100	1.05	4.30	1.22	9.74	33.6285	18.04

**Figure 5:** Ensemble averages vs time of MHAAR-RB, MwG, and AIS MCMC. Left: All the settings, Right: Best settings

justifies the use of both annealing (via AIS MCMC) and averaging (via MHAAR-RB), especially the latter proves more useful owing to the relative ease of implementing the averaging compared to annealing. Also, we provide the averages for the two parameters where the difference is most visible; for the other parameters the algorithms showed similar performance.

The other details of our experiment are as follows. The model parameters are selected in parallel with Neal [2010]: We take  $\varsigma = 0.01$ ,  $v = 1$ , and the prior distribution the vector  $\log \theta$  is taken a normal distribution with mean  $\log 0.5$  and unit variance for each component, with a correlation of 0.69 for any pair of components. The other parameters  $\tau, \sigma^2$  are apriori independent from  $\theta$  and among themselves, with  $\log \tau \sim \mathcal{N}(0, 2.25)$  and  $\log \sigma \sim \mathcal{N}(\log 0.5, 2.25)$ . AIS MCMC and MHAAR-RB attempt to update one component of  $\theta$  at a time with the same proposal mechanism. For each component, a normal random walk proposal is used for  $\log \theta_i$  with mean 0 and standard deviation 2. At the intermediate steps of AIS MCMC, the fast variables are updated with an MH kernel with random walk proposals on  $\log \tau$  and  $\log \sigma$  with zero mean and standard deviations 0.6 for both. We run MHAAR-RB with no annealing, i.e.,  $\gamma_{\theta, \vartheta} = \gamma_{\theta}$ , ending up with the acceptance ratio in (23) with  $T = 1$  (Superiority of no annealing in general is shown in the previous example). Moreover,  $q_{\theta, \vartheta}(z)$  is taken as density of the prior distribution of  $z$ , therefore, we have  $Q_1^M = Q_2^M$ .

## 4 State-space models: SMC and cSMC within MHAAR

In Sections 2 and 3, we have shown how two different generic MHAAR strategies which consist of averaging estimates of the acceptance ratio could be helpful. Here we extend the methodology in Section 3 to state-space models. Specifically we present methods where dependent acceptance ratios arising from a single conditional SMC algorithm can be averaged in order to improve performance.

## 4.1 State-space models and cSMC

In its simplest form, a state-space model (SSM) is comprised of a latent Markov chain  $\{Z_t; t \geq 1\}$  taking its values in some measurable space  $(Z, \mathcal{Z})$  and observations  $\{Y_t; t \geq 1\}$  taking values in  $(Y, \mathcal{Y})$ . The latent process has initial probability of density  $f_\theta(z_1)$  and transition density  $f_\theta(z_{t-1}, z_t)$ , dependent on a parameter  $\theta \in \Theta \subset \mathbb{R}^{d_\theta}$ . An observation at time  $t \geq 1$  is assumed conditionally independent of all other random variables given  $Z_t = z_t$  and its conditional observation density is  $g_\theta(z_t, y_t)$ . The corresponding joint density of the latent and observed variables up to time  $T \geq 1$  is

$$p_\theta(z_{1:T}, y_{1:T}) = f_\theta(z_1) \prod_{t=2}^T f_\theta(z_{t-1}, z_t) \prod_{t=1}^T g_\theta(z_t, y_t). \quad (30)$$

The densities  $f_\theta$  and  $g_\theta$  could also depend on  $t$ , at the expense of notational complications. In order to alleviate notation and ensure consistency we let  $z := z_{1:T}$  and  $y := y_{1:T}$ . The likelihood function associated to the observations  $y$  can be obtained

$$\ell_\theta(y) := \int_{Z^T} p_\theta(z, y) dz. \quad (31)$$

With a prior  $\eta(d\theta)$  on  $\theta$  with density  $\eta(\theta)$ , the joint posterior  $\pi(d(\theta, z))$  has the density

$$\pi(\theta, z) \propto \eta(\theta) p_\theta(z, y)$$

so that  $\pi(\theta) \propto \eta(\theta) \ell_\theta(y)$  and  $\pi_\theta(z) := p_\theta(z | y) = p_\theta(z, y) / \ell_\theta(y)$ . Therefore, the acceptance ratio of the marginal MCMC algorithm for SSM can be written as

$$r(\theta, \vartheta) = \frac{q(\vartheta, \theta) \eta(\vartheta) \ell_\vartheta(y)}{q(\theta, \vartheta) \eta(\theta) \ell_\theta(y)}. \quad (32)$$

Conditional sequential Monte Carlo (cSMC) introduced in [Andrieu et al. \[2010\]](#) is an MCMC transition kernel akin to particle filters that is particularly well suited to sampling from  $\pi_\theta(dz)$ . It was shown in [Lindsten and Schön \[2012\]](#) that cSMC with backward sampling [[Whiteley, 2010](#)] can be used efficiently as part of a more elaborate Metropolis-within-Particle Gibbs algorithm in order to sample from the posterior distribution  $\pi(d(\theta, z))$ ; see Algorithm 4. In [Gunawan et al. \[2020\]](#) it is shown how this can be combined with ideas of [Deligiannidis et al. \[2018\]](#) to improve performance in specific scenarios.

The cSMC algorithm with backward sampling for state-space models used to present our results is given in Algorithm 6 in Appendix C. To simplify exposition we consider the bootstrap particle filter where the particles are initialised according to  $f_\theta(z_1)$  and propagated according to the state transition  $f_\theta(z_{t-1}, z_t)$ ; our results can be extended straightforwardly to other choices. The cSMC produces  $T \times M$  samples from which  $M^T$  paths can be sampled using the backward recursion of [[Whiteley, 2010](#)]. The cSMC returns only one such path when used in Algorithm 4, which may seem to be wasteful. A natural idea is to make use of multiple—or even all  $M^T$  possible—trajectories and average out the corresponding acceptance ratios (33) before accepting or rejecting. We show that this is indeed possible theoretically with Algorithms 5 and 7 in the next section. We then

---

**Algorithm 4:** Metropolis-within-particle Gibbs

---

**Input:** Current sample  $(\theta, z)$

**Output:** New sample

- 1 Sample  $\vartheta \sim q(\theta, \cdot)$  and  $z' \sim \text{cSMC}(M, \theta, z)$ .
- 2 Return  $(\vartheta, z')$  with probability

$$\min \left\{ 1, \frac{q(\vartheta, \theta) \eta(\vartheta) p_{\vartheta}(z', y)}{q(\theta, \vartheta) \eta(\theta) p_{\theta}(z', y)} \right\}; \quad (33)$$

otherwise return  $(\theta, z')$ .

---

show that these schemes are highly advantageous on parallel computing architecture, but also on serial machines in some difficult scenarios. The justification of the algorithms is postponed to Appendix C; while this can be thought of as extensions of the results of Section 3 the dependence structure implied by the cSMC leads to significant conceptual and notational complications.

## 4.2 MHAAR with cSMC for state-space models

We will first present an unbiased estimator of the marginal acceptance ratio in (32) for SSM using particles produced by a cSMC iteration. Building on this we present our MHAAR algorithm for SSM.

### 4.2.1 Unbiased estimator of the acceptance ratio using particles of cSMC

The particles  $\mathbf{v} = v_{1:T}^{(1:M)}$  outputted by the cSMC update can be partitioned as  $\mathbf{v} = (z, \mathbf{u})$ , where  $z := v^{(1)}$  is the path conditional upon which the cSMC is run, and  $\mathbf{u} := v^{(\bar{1})}$  consists of the rest of the variables in  $\mathbf{v}$ . It can be shown that the conditional distribution of  $\mathbf{u}$  given  $(\theta, z) \in \Theta \times \mathcal{Z}$  is given by

$$\Phi_{\theta}(z, d\mathbf{u}) = \prod_{i=2}^M f_{\theta}(dv_1^{(i)}) \prod_{t=2}^T \left\{ \prod_{i=2}^M \frac{\sum_{j=1}^M w_{t-1, \theta}(v_{t-1}^{(j)}) f_{\theta}(v_{t-1}^{(j)}, dv_t^{(i)})}{\sum_{j=1}^M w_{t-1, \theta}(v_{t-1}^{(j)})} \right\}.$$

The law of the indices  $\mathbf{k} := (k_1, \dots, k_T)$  drawn in the backward sampling step in Algorithm 6 (lines 10-14) conditional upon  $\theta$  and  $\mathbf{v}$  is given by

$$b_{\theta}(\mathbf{k}|\mathbf{v}) := \frac{w_{T, \theta}(v_T^{(k_T)})}{\sum_{i=1}^M w_{T, \theta}(v_T^{(i)})} \prod_{t=1}^{T-1} \frac{w_{t, \theta}(v_t^{(k_t)}) f_{\theta}(v_t^{(k_t)}, v_{t+1}^{(k_{t+1})})}{\sum_{i=1}^M w_{t, \theta}(v_t^{(i)}) f_{\theta}(v_t^{(i)}, v_{t+1}^{(k_{t+1})})}.$$

Further, for any  $\theta, \vartheta, \zeta \in \Theta$ , and  $z, z' \in \mathcal{Z}^T$ , define

$$r_{z, z'}(\theta, \vartheta; \zeta) = \frac{q(\vartheta, \theta) \eta(\vartheta) p_{\vartheta}(z', y) p_{\zeta}(z, y)}{q(\theta, \vartheta) \eta(\theta) p_{\theta}(z', y) p_{\theta}(z, y)}. \quad (34)$$

In the following, we show that it is possible to construct unbiased estimators of  $r(\theta, \vartheta)$  in (32) using cSMC, provided we have a random sample  $z \sim \pi_{\theta}(\cdot)$ . Specifically, this is obtained as the expected value of  $r_{z, v(\mathbf{k})}(\theta, \vartheta; \zeta)$  with respect to the backward sampling distribution of  $\mathbf{k}$ ,  $b_{\theta}(\mathbf{k}|\mathbf{v})$ .

**Theorem 4.** For  $\theta, \vartheta, \zeta \in \Theta$  and any  $M \geq 1$ , let  $z \sim \pi_\theta(\cdot)$ ,  $\mathbf{v}|z \sim \text{cSMC}(M, \zeta, z)$  be the generated particles from the cSMC algorithm targeting  $\pi_\zeta(\cdot)$  with  $M$  particles, conditioned on  $z$ . Then,  $r_{\mathbf{1}, \mathbf{v}}(\theta, \vartheta; \zeta)$  is an unbiased estimator of  $r(\theta, \vartheta)$  in (32), where for  $\mathbf{v} \in \mathbb{Z}^{TM}$ ,  $\mathbf{l} \in \llbracket M \rrbracket^T$ ,  $r_{\mathbf{l}, \mathbf{v}}(\theta, \vartheta; \zeta)$  is defined as

$$r_{\mathbf{l}, \mathbf{v}}(\theta, \vartheta; \zeta) := \sum_{\mathbf{k} \in \llbracket M \rrbracket^T} r_{v^{(1)}, v^{(\mathbf{k})}}(\theta, \vartheta; \zeta) b_\zeta(\mathbf{k}|\mathbf{v}). \quad (35)$$

The proof of Theorem 4 is left to Appendix C.1. Theorem 4 is original to the best of our knowledge and we find it interesting in several aspects. Firstly, unlike the estimator in Metropolis-within-Particle Gibbs (Algorithm 4), the estimator in Theorem 4 uses all possible paths from the particles generated by the cSMC. Also, with a slight modification one can similarly obtain unbiased estimators for  $\pi(\vartheta)/\pi(\theta)$  which is of primary interest in some applications. The theorem is derived from Del Moral et al. [2010, Theorem 5.2] and the results in Andrieu et al. [2010] relating the laws of cSMC and SMC.

#### 4.2.2 MHAAR-RB for SSM

Theorem 4 motivates the design of a MHAAR algorithm using the unbiased estimator (35) as its acceptance ratios. We describe the algorithm, MHAAR-RB for SSM, in detail below. The procedure requires a pair of functions  $\zeta_1 : \Theta^2 \rightarrow \Theta$  and  $\zeta_2 : \Theta^2 \rightarrow \Theta$  satisfying  $\zeta_1(\theta, \vartheta) = \zeta_2(\vartheta, \theta)$  for  $\theta, \vartheta \in \Theta$ , in order to determine the intermediate parameter value for which the cSMC is run. MHAAR-RB for SSM targets the joint distribution for the variable  $\xi = (\theta, \vartheta, \mathbf{v}, \mathbf{k}, c) \in \Theta^2 \times \mathbb{Z}^{MT} \times \llbracket M \rrbracket^T \times \{1, 2\}$  defined as

$$\hat{\pi}(\mathrm{d}(\theta, \vartheta, \mathbf{v}, \mathbf{k}, c)) = \frac{1}{2} \pi(\mathrm{d}(\theta, z)) Q_c^M(\theta, z; \mathrm{d}(\mathbf{u}, \mathbf{k})). \quad (36)$$

where we have used  $z := v^{(1)}$ . Clearly, the marginal distribution for  $(\theta, z)$  is  $\pi(\mathrm{d}(\theta, z))$ , as desired. The proposal mechanisms are

$$Q_c^M(\theta, z; \mathrm{d}(\mathbf{u}, \mathbf{k})) = q(\theta, \mathrm{d}\vartheta) \Phi_{\zeta_c(\theta, \vartheta)}(z, \mathrm{d}\mathbf{u}) b_{\theta, \vartheta}^{(c)}(\mathbf{k}|\mathbf{v}), \quad c = 1, 2,$$

where the sampling probabilities are given as  $b_{\theta, \vartheta}^{(2)}(\mathbf{k}|\mathbf{v}) = b_{\zeta_2(\theta, \vartheta)}(\mathbf{k}|\mathbf{v})$  and

$$b_{\theta, \vartheta}^{(1)}(\mathbf{k}|\mathbf{v}) = \frac{r_{v^{(1)}, v^{(\mathbf{k})}}(\theta, \vartheta; \zeta_1(\theta, \vartheta)) b_{\zeta_1(\theta, \vartheta)}(\mathbf{k}|\mathbf{v})}{r_{\mathbf{1}, \mathbf{v}}(\theta, \vartheta; \zeta_1(\theta, \vartheta))},$$

which is obtained by weighting the backward sampling probabilities of the cSMC by the acceptance ratios they correspond to, yielding the normalising constant  $r_{\mathbf{1}, \mathbf{v}}(\theta, \vartheta; \zeta_1(\theta, \vartheta))$  defined in (35). One iteration of MHAAR-RB for SSM consists of the following main steps:

1. Sample  $c \sim \text{Unif}(\{1, 2\})$ , then sample  $(\vartheta, \mathbf{u}, \mathbf{k}) \sim Q_c^N(\theta, z; \cdot)$ , and form  $\xi = (\theta, \vartheta, \mathbf{v}, \mathbf{k}, c)$ .
2. Propose an MH update of  $\xi$  via the involution

$$\xi' = \varphi(\theta, \vartheta, \mathbf{v}, \mathbf{k}, c) := (\vartheta, \theta, \mathbf{s}_{\mathbf{1}, \mathbf{k}}(\mathbf{v}), \mathbf{k}, 3 - c), \quad (37)$$

where  $\mathbf{s}_{\mathbf{1}, \mathbf{k}}(\mathbf{v})$  is an operator on  $\mathbf{v}$  that swaps  $v^{(1)}$  and  $v^{(\mathbf{k})}$ .

3. Accept  $\xi'$  with acceptance probability  $\min\{1, \hat{r}(\xi)\}$ , otherwise reject and keep  $\xi$ .

We prove in Appendix C.2.1 that this proposed involution leads to the averaged acceptance ratio  $r_{\mathbf{1}, \mathbf{v}}(\theta, \vartheta; \zeta_1(\theta, \vartheta))$  in its acceptance probabilities, as stated in the theorem below.

**Theorem 5.** *With the joint distribution  $\hat{\pi}$  defined in (36), the acceptance ratio for the proposed involution defined in (37) is given by*

$$\hat{r}(\xi) := \begin{cases} r_{\mathbf{1}, \mathbf{v}}(\theta, \vartheta; \zeta_1(\theta, \vartheta)), & c = 1, \\ 1/r_{\mathbf{k}, \mathbf{v}}(\vartheta, \theta; \zeta_2(\theta, \vartheta)), & c = 2. \end{cases}$$

The proof has two interesting by-products: (i) An alternative proof of Theorem 4, and (ii) another unbiased estimator of  $r(\theta, \vartheta)$  which uses all  $M^T$  possible paths formed from the particles generated by the cSMC, which is we state precisely in the following corollary.

**Corollary 1.** *For  $\theta, \vartheta, \zeta \in \Theta$  and any  $M \geq 1$ , let  $z \sim \pi_\theta(\cdot)$ ,  $\mathbf{v}|z \sim \text{cSMC}(M, \zeta, z)$  be the generated particles from the cSMC algorithm with  $M$  particles conditional on  $\zeta, z$  and  $\mathbf{k}|\mathbf{v} \sim b_\zeta(\cdot|\mathbf{v})$ . Then,  $1/r_{\mathbf{k}, \mathbf{v}}(\vartheta, \theta; \zeta)$  is an unbiased estimator of  $r(\theta, \vartheta)$ .*

We present MHAAR-RB for SSM in Algorithm 5. The per iteration computational complexity of Algorithm 5 is  $\mathcal{O}(M^2T)$ . This follows upon observing that the unnormalised probability in (35) can be written as

$$r_{v(1), v(\mathbf{k})}(\theta, \vartheta; \zeta) b_\zeta(\mathbf{k}|\mathbf{v}) =: \varrho_{\mathbf{1}, \mathbf{v}}(\mathbf{k}) = \varrho_{\mathbf{1}, \mathbf{v}, 1}(k_1) \prod_{t=2}^T \varrho_{\mathbf{1}, \mathbf{v}, t}(k_{t-1}, k_t)$$

for an appropriate choice of the functions  $\varrho_{\mathbf{1}, \mathbf{v}, t}$  and that  $r_{\mathbf{1}, \mathbf{v}}(\theta, \vartheta; \zeta) = \sum_{\mathbf{k} \in [M]^T} \varrho_{\mathbf{1}, \mathbf{v}}(\mathbf{k})$  can be computed using a sum-product algorithm, while sampling  $\mathbf{k}$  with probability proportional to  $\varrho_{\mathbf{1}, \mathbf{v}}(\mathbf{k})$ , required when  $c = 1$ , can be performed with a forward-filtering backward-sampling algorithm [Zucchini et al., 2016]. We note that: (a) while complexity is  $\mathcal{O}(M^2T)$  the operations involved are often much cheaper than for the cSMC since, for example, likelihood terms involved need not re-evaluation, (b) recent work investigates the implementation of such recursions on GPUs e.g. Natarajan and Chandrachoodan [2018], although this is far beyond the scope of the present methodological paper.

**Refreshing  $z$  via delayed rejection:** In Section 3, in the particular scenario where the latent variable sequence consists of iid states, we have already discussed how a delayed rejection step can be included to refresh the variable  $z$  upon a ‘stage 1’ rejection, at a minimal computational cost. Delayed rejection is also possible for SSM and is particularly attractive when  $c = 1$  and  $\zeta_1(\theta, \vartheta) = \theta$ . In this case  $z$  can be refreshed upon rejection by simply performing another backward sampling iteration on the already sampled particles  $\mathbf{v}$ . Otherwise a second accept/reject step is required. The proof of validity for all scenarios is left to Appendix C.2.2. The delayed rejection step is included in Algorithm 5 as an ‘optional’ step and its cost is  $\mathcal{O}(MT)$ .

---

**Algorithm 5:** MHAAR-RB for SSM

---

**Input:** Current sample  $(\theta, z)$

**Output:** New sample

- 1 Sample  $\vartheta \sim q(\theta, \cdot)$  and  $c \sim \text{Unif}(\{1, 2\})$ , and set  $\zeta = \zeta_c(\theta, \vartheta)$ .
  - 2 **if**  $c = 1$  **then**
  - 3     Run a cSMC( $M, \zeta, z$ ) targeting  $\pi_\zeta$  conditional on  $z$  to obtain  $\mathbf{v}$ .
  - 4     Sample  $\mathbf{k} \sim b_{\theta, \vartheta}^{(1)}(\cdot | \mathbf{v})$  and set  $z' = v^{(\mathbf{k})}$
  - 5     Return  $(\vartheta, z')$  with probability  $\min\{1, r_{1, \mathbf{v}}(\theta, \vartheta; \zeta)\}$ ; otherwise return  $(\theta, z)$ .
  - 6 **else**
  - 7     Run a cSMC( $M, \zeta, z$ ) targeting  $\pi_\zeta$  conditional upon  $z$  to obtain  $\mathbf{v}$ .
  - 8     Sample  $\mathbf{k} \sim b_{\theta, \vartheta}^{(2)}(\cdot | \mathbf{v})$  and set  $z' = v^{(\mathbf{k})}$ .
  - 9     Return  $(\vartheta, z')$  with probability  $\min\{1, 1/r_{\mathbf{k}, \mathbf{v}}(\vartheta, \theta; \zeta)\}$ ; otherwise return  $(\vartheta, z')$ .
  - 10 **Optional refreshment of**  $z$
  - 11 **if** *the move is rejected*,  $c = 1$ , *and*  $\zeta = \theta$  **then**
  - 12     Sample  $\mathbf{l} \sim b_\theta(\cdot | \mathbf{v})$  and return  $(\theta, v^{(\mathbf{l})})$ .
- 

**Example 10.** We consider the following linear Gaussian SSM

$$\begin{aligned} Z_t &= \phi(Z_{t-1} - (1 - a)\theta) + (1 - a)\theta + V_t, \quad t \geq 2 \\ Y_t &= Z_t + a\theta + W_t, \quad t \geq 1. \end{aligned}$$

where  $\phi > 0$  is a coefficient,  $a \in [0, 1]$ ,  $Z_1 \sim \mathcal{N}(0, \sigma_z^2)$ ,  $V_t \stackrel{\text{iid}}{\sim} \mathcal{N}(0, (1 - \phi^2)\sigma_z^2)$ , and  $W_t \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma_y^2)$ . Naturally a Kalman filter can be used here to compute the likelihood function efficiently and no Monte Carlo methods are needed. However this model offers a fully controllable testbed useful to illustrate the type of situations where MHAAR is of interest. Importantly the likelihood function  $\theta \mapsto \ell_\theta(y, a)$  does not depend on the choice of  $a$  but, assuming a prior distribution on  $\theta$ , the posterior dependency between  $\theta$  and  $Z_{1:T}$  does. As a result the mixing properties of a Gibbs sampler sampling alternately from  $\pi(\theta | z_{1:T})$  and  $\pi(z_{1:T} | \theta)$  are highly dependent on the choice of  $a$ . For example for  $\phi = 0$ , Papaspiliopoulos et al. [2003] showed that for  $\sigma_z^2/\sigma_y^2 \ll 1$  (resp.  $\sigma_z^2/\sigma_y^2 \gg 1$ ) the choice  $a \approx 1$  ( $a \approx 0$ ) leads to strong posterior dependence.

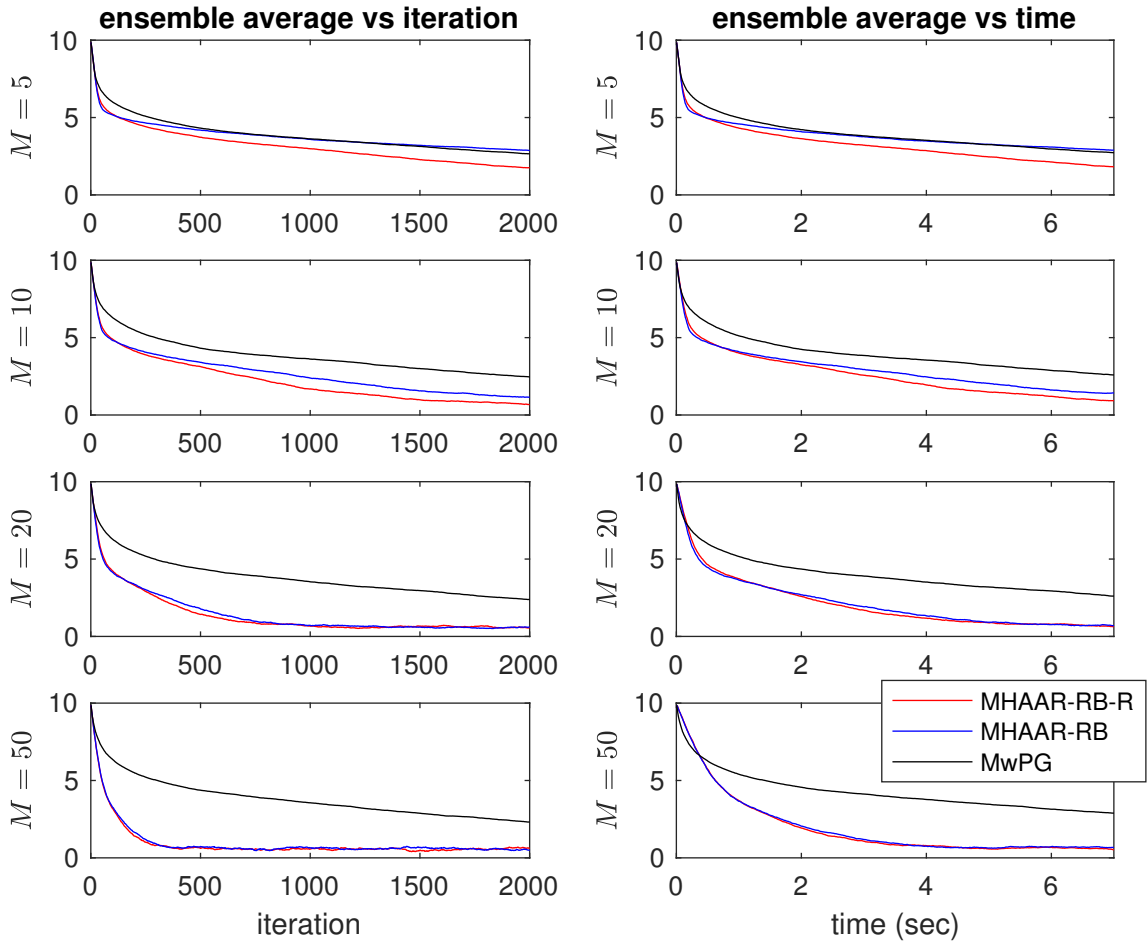
We generated a dataset of size  $T = 100$  from this SSM with  $\theta^* = 1$ ,  $\phi = 0.95$  and noise parameters  $\sigma_z^2 = 1$  and  $\sigma_y^2 = 0.1$ , the regime where  $a \approx 1$  leads to strong dependence, hence our choice of  $a = 1$ . We compared MHAAR-RB, MHAAR-RB-R (with refreshment of  $z$  upon rejection) for SSM as in Algorithm 5 and MwPG in terms of IAC time and IAC  $\times$  CPU time per iteration for  $\theta$  for different values of  $M$ . For MHAAR-RB and MHAAR-RB-R, we used  $\zeta_1(\theta, \vartheta) = \theta$ . Each run is performed for  $10^6$  iterations, except that we run MwPG for  $5 \times 10^6$  iterations to overcome the variability in the estimates for the IAC times. The prior for  $\theta$  is taken as  $\mathcal{N}(0, 10^4)$ . For all the algorithms, a random walk proposal is used with a proposal standard deviation of  $\sigma_q = 0.3$ . The results are displayed in Table 3. We observe that MHAAR-RB and MHAAR-RB-R's response to increasing  $M$  is substantial and should be contrasted with the standard MwPG's underwhelming performance. Further we note the superiority of MHAAR-RB and MHAAR-RB-R on

MwPG even when the IAC time is rescaled with the computation time, that is MHAAR-RB and MHAAR-RB-R outperform MwPG even on a serial machine for this example.

Figure 6 shows the ensemble averages over 100 runs (see Example 7) for the posterior expectation of  $\theta$  versus iteration number and time for the three algorithms, illustrating burn-in length. The results mirror those of Table 3 concerned with IAC times with MHAAR-RB and MHAAR-RB-R vastly superior to MwPG in terms of burn in length, with much better reactivity to increasing  $M$ .

$M$	IAC time ( $\times 10^3$ )			IAC $\times$ CPU time per iteration		
	MHAAR-RB	MHAAR-RB-R	MwPG	MHAAR-RB	MHAAR-RB-R	MwPG
5	4.1801	1.7666	4.2519	13.8732	5.7242	13.2968
10	1.4971	1.1556	3.8536	5.8461	4.4266	12.6457
20	0.4713	0.4332	3.5337	2.7598	2.5566	12.1545
50	0.1579	0.1516	3.2501	1.7587	1.7935	14.0562

**Table 3:** Comparison of MHAAR-RB, MHAAR-RB-R, and MwPG in terms of IAC and IAC  $\times$  CPU time per iteration for  $\theta$ .



**Figure 6:** Ensemble averages for the posterior expectation of  $\theta$  vs iteration number and time for the algorithms compared in Table 3.



### 4.2.3 Reduced computational cost via subsampling

The  $\mathcal{O}(M^2T)$  cost per iteration of MHAAR-RB for SSM precludes its application as  $M$  becomes large, as required in some applications. A computationally less demanding and intuitive version of Algorithm 5 could use a subsampled version of the large sum in (35) applying the backward sampling procedure  $N$  times to recover  $N$  paths. That is, letting  $\mathbf{u} = (u^{(1)}, \dots, u^{(N)}) \in \mathbb{Z}^{TN}$ , a natural idea is to use the unbiased estimator of (35)

$$r_{z,\mathbf{u}}^N(\theta, \vartheta; \zeta) = \frac{1}{N} \sum_{i=1}^N r_{z,u^{(i)}}(\theta, \vartheta; \zeta), \quad (38)$$

where

$$u^{(1)}, \dots, u^{(N)} \stackrel{\text{iid}}{\sim} \sum_{\mathbf{k} \in \llbracket M \rrbracket^T} b_{\zeta}(\mathbf{k}|\mathbf{v}) \delta_{v(\mathbf{k})}(\cdot).$$

Designing an algorithm using this acceptance ratio (38) while preserving the correct invariant distribution  $\pi(d(\theta, z))$  is possible in the MHAAR framework. The resulting algorithm, which we name MHAAR-S(ubsample) for SSM, is presented in Algorithm 7 in Appendix C.3. The computational complexity of MHAAR-S for SSM is  $\mathcal{O}(NMT)$  per iteration instead of  $\mathcal{O}(M^2T)$  for Algorithm 7. We note again that sampling  $N$  paths using backward sampling is an embarrassingly parallelisable operation. Details and correctness of MHAAR-S for SSM as well as additional numerical results are provided in Appendix C.3.1.

**Example 11 (Example 10, ctd).** We run MHAAR-S for HMM for the dataset used in Example 10 with  $M = 20$  particles and several values of  $N$ . Table 4 shows the IAC times for MHAAR-S for SSM, estimated from  $2 \times 10^6$  iterations, in comparison with IAC times of MHAAR-RB-R and MwPG with the same number of particles. We also show the ensemble averages of those algorithms, obtained from 100 independent runs, in Figure 7. Note that, using all the  $M^T$  possible paths, the MHAAR-RB and MHAAR-RB-R algorithms set a limit on the performance of MHAAR-S for SSM. Both the table and the figure show that using multiple paths results in gains in terms of convergence to equilibrium compared to MwPG, illustrating the potential of the MAHHR approach to leverage massively parallel architectures and reduce wall-clock time.

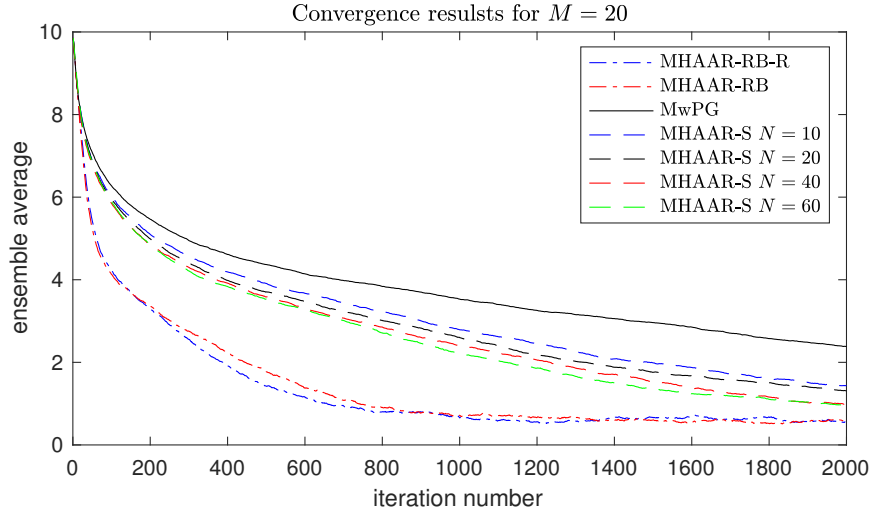
MHAAR-S				MHAAR-RB-R	MHAAR-RB	MwPG
$N = 10$	$N = 20$	$N = 40$	$N = 60$			
2.0378	1.5770	1.5507	1.4047	0.4332	0.4713	3.5337

**Table 4:** Comparison of MHAAR-S and MwPG in terms of IAC time ( $\times 10^3$ ). Each run is performed for 500000 iterations.  $M = 20$  is taken for all runs.

## 5 Discussion

In this paper, we exploit the ability to use more than one proposal schemes within a MH update. We derive several useful MHAAR algorithms that enable averaging multiple





**Figure 7:** Ensemble averages for the posterior expectation of  $\theta$  vs iteration number and time for MHAAR-S, in comparison with MHAAR-RB and MwPG.

estimates of acceptance ratios, which would not be valid by using a standard single proposal MH update. The framework of MHAAR is rather general and provides a generic way of improving performance of MH update based algorithm for a wide range of problems. This is illustrated with doubly intractable models, general latent variable models, trans-dimensional models, and general state-space models. Although relevant in specific scenarios involving computations on serial machines, MHAAR algorithms are particularly useful when implemented on a parallel architecture since the computation required to have an average acceptance ratio estimate can largely be parallelised. In particular our experiments demonstrate significant reduction of the burn in period required to reach equilibrium, an issue for which very few generic approaches exist currently.

## 6 Acknowledgements

CA and SY acknowledge support from EPSRC “Intractable Likelihood: New Challenges from Modern Applications (ILike)” (EP/K014463/1) and the Isaac Newton Institute for Mathematical Sciences, Cambridge, for support and hospitality during the programme “Scalable inference; statistical, algorithmic, computational aspects” during which some the work was carried out (EPSRC grant EP/K032208/1). CA and AD acknowledge support of EPSRC grants Bayes4Health (EP/R018561/1) and CoSInES (EP/R034710/1). NC is partially supported by a grant from the French National Research Agency (ANR) as part of program ANR-11-LABEX-0047. The authors would also like to thank Nick Whiteley for useful discussions.

## References

Christophe Andrieu and Gareth O. Roberts. The pseudo-marginal approach for efficient Monte Carlo computations. *Annals of Statistics*, 37(2):569–1078, 2009.

- Christophe Andrieu and Matti Vihola. Establishing some order amongst exact approximations of MCMCs. *Annals of Applied Probability*, 26(5):2661–2696, 10 2016.
- Christophe Andrieu, Arnaud Doucet, and Roman Holenstein. Particle Markov chain Monte Carlo methods. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 72:269–342, 2010. doi: 10.1111/j.1467-9868.2009.00736.x.
- Christophe Andrieu, Arnaud Doucet, Sinan Yıldırım, and Nicolas Chopin. On the utility of Metropolis-Hastings with asymmetric acceptance ratio. *ArXiv e-prints*, (1803.09527), 2018.
- Christophe Andrieu, Anthony Lee, and Sam Livingstone. A general perspective on the Metropolis-Hastings kernel. *ArXiv e-prints*, 2020.
- M. Beaumont. Estimation of population growth of decline in genetically monitored populations. *Genetics*, 164:1139–1160, 2003.
- Mark A. Beaumont, Wenyang Zhang, and David J. Balding. Approximate Bayesian computation in population genetics. *Genetics*, 162(4):2025–2035, December 2002. ISSN 0016-6731. URL <http://www.genetics.org/content/162/4/2025.abstract>.
- Luke Bornn, Natesh S Pillai, Aaron Smith, and Dawn Woodard. The use of a single pseudo-sample in approximate Bayesian computation. *Statistics and Computing*, 27(3):583–590, 2017.
- Joe Caaney. *Contributions to Exact Approximation Methodology*. PhD thesis, University of Bristol, School of Mathematics, University of Bristol, 2013.
- Pierre Del Moral, Arnaud Doucet, and Sumeetpal S Singh. A backward particle interpretation of Feynman-Kac formulae. *ESAIM: Mathematical Modelling and Numerical Analysis*, 44(5):947–975, 2010.
- Georgios Deligiannidis, Arnaud Doucet, and Michael K. Pitt. The correlated pseudo-marginal method. *Journal of the Royal Statistical Society. Series B: Statistical Methodology*, 7 2018. ISSN 1369-7412. doi: 10.1111/rssb.12280.
- Petros Dellaportas and Ioannis Kontoyiannis. Control variates for estimation based on reversible Markov chain Monte Carlo samplers. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 74(1):133–161, 2012.
- Jean-François Delmas and Benjamin Jourdain. Does waste recycling really improve the multi-proposal Metropolis-Hastings algorithm? An analysis based on control variates. *Journal of Applied Probability*, 46(4):938–959, 2009.
- P. Green. Reversible jump Markov chain Monte Carlo for Bayesian model determination. *Biometrika*, 82(4):711–732, 1995.
- David Gunawan, Chris Carter, and Robert Kohn. On scalable particle Markov chain Monte Carlo, 2020.

- G. Karagiannis and C. Andrieu. Annealed importance sampling for reversible jump MCMC algorithms. *Journal of Computational and Graphical Statistics*, 22(3):623–648, 2013.
- Anthony Lee, Christopher Yau, Michael B Giles, Arnaud Doucet, and Christopher C Holmes. On the utility of graphics cards to perform massively parallel simulation of advanced Monte Carlo methods. *Journal of computational and graphical statistics*, 19(4):769–789, 2010.
- David A Levin and Yuval Peres. *Markov chains and mixing times*, volume 107. American Mathematical Soc., 2017.
- F. Lindsten and T. B. Schön. On the use of backward simulation in the particle Gibbs sampler. In *2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 3845–3848, March 2012. doi: 10.1109/ICASSP.2012.6288756.
- Paul Marjoram, John Molitor, Vincent Plagnol, and Simon Tavaré. Markov chain Monte Carlo without likelihoods. *Proceedings of the National Academy of Sciences of the United States of America*, 100(26):15324–15328, 2003. ISSN 00278424. URL <http://www.jstor.org/stable/3149004>.
- J. Møller, A. N. Pettitt, R. Reeves, and K. K. Berthelsen. An efficient Markov chain Monte Carlo method for distributions with intractable normalising constants. *Biometrika*, 93(2):451–458, 2006. doi: 10.1093/biomet/93.2.451. URL <http://biomet.oxfordjournals.org/content/93/2/451.abstract>.
- A Müller and D Stoyan. Comparison methods for stochastic models and risks. *John Wiley&Sons Ltd., Chichester*, 2002.
- I. Murray, Z. Ghahramani, and D. J. C. MacKay. MCMC for doubly-intractable distributions. In *Proceedings of the 22nd Annual Conference on Uncertainty in Artificial Intelligence (UAI-06)*, pages 359–366, 2006.
- K. Natarajan and N. Chandrachoodan. Lossless parallel implementation of a turbo decoder on GPU. In *2018 IEEE 25th International Conference on High Performance Computing (HiPC)*, pages 133–142, 2018. doi: 10.1109/HiPC.2018.00023.
- Radford M. Neal. Taking bigger Metropolis steps by dragging fast variables. Technical report, University of Toronto, 2004.
- Radford M. Neal. MCMC using ensembles of states for problems with fast and slow variables such as Gaussian process regression. Technical report, University of Toronto, 2010.
- O. Papaspiliopoulos, G.O. Roberts, and M. Skold. Non-centred parameterisations for hierarchical models and data augmentation. In J. M. Bernardo, M. J. Bayarri, J. O. Berger, A. P. Dawid, D. Heckerman, A. F. M. Smith, and M. West, editors, *Bayesian Statistics VII*, pages 307–327. 2003.

- J. Pritchard, M. Seielstad, A. Perez-Lezaun, and M. Feldman. Population growth of human Y chromosomes: a study of Y chromosome microsatellites. *Molecular Biology and Evolution*, 16:1791–1798, 1999.
- Chris Sherlock, Alexandre H. Thiery, and Anthony Lee. Pseudo-marginal Metropolis-Hastings sampling using averages of unbiased estimators. *Biometrika*, 104(3):727–734, 2017. doi: 10.1093/biomet/asx031. URL [+http://dx.doi.org/10.1093/biomet/asx031](http://dx.doi.org/10.1093/biomet/asx031).
- Andrew Sohn. Parallel n-ary speculative computation of simulated annealing. *IEEE Transactions on Parallel and Distributed systems*, 6(10):997–1005, 1995.
- Marc A Suchard, Quanli Wang, Cliburn Chan, Jacob Frelinger, Andrew Cron, and Mike West. Understanding GPU programming for statistical computation: Studies in massively parallel massive mixtures. *Journal of computational and graphical statistics*, 19(2):419–438, 2010.
- Luke Tierney. A note on Metropolis Hastings kernels for general state spaces. *Annals of Applied Probability*, 8(1):1–9, 1998.
- Nick Whiteley. Discussion on particle Markov chain Monte Carlo methods. *Journal of the Royal Statistical Society: Series B*, 72(3):306–307, 2010.
- Sinan Yildirim, Sumeetpal S. Singh, Thomas Dean, and Ajay Jasra. Parameter estimation in hidden Markov models with intractable likelihoods using sequential Monte Carlo. *Journal of Computational and Graphical Statistics*, 24(3):846–865, 2015. doi: 10.1080/10618600.2014.938811. URL <https://doi.org/10.1080/10618600.2014.938811>.
- Sinan Yildirim, Christophe Andrieu, and Arnaud Doucet. Scalable Monte Carlo inference for state-space models, 2018.
- Giacomo Zanella. Informed proposals for local MCMC in discrete spaces. *Journal of the American Statistical Association*, 115(530):852–865, 2020. doi: 10.1080/01621459.2019.1585255. URL <https://doi.org/10.1080/01621459.2019.1585255>.
- W. Zucchini, I. MacDonald, and R. Langrock. *Hidden Markov models for time series: An introduction using R*, volume 6. Chapman and Hall/CRC, 2016. doi: <https://doi.org/10.1201/b20790>.

## A Proofs for the theorems in Section 2

### A.1 Acceptance ratio of Algorithm 1

*Proof of Theorem 1.* Let  $\tilde{\pi}_0(d(\theta, z, u)) = \pi(d(\theta, z))q(\theta, d\vartheta)Q_{\theta, \vartheta, z}(du)$ . Then,  $r_{u^{(i)}}(\theta, \vartheta, z)$  is the acceptance ratio for  $\tilde{\pi}_0$  corresponding to the involution  $\varphi_0(\theta, \vartheta, z, u) = (\vartheta, \theta, \phi_{\theta, \vartheta}(z, u))$ .

Also, observe that, when  $c = 1$ , we have

$$\begin{aligned}\pi(d\xi) &= \hat{\pi}_0(d(\theta, \vartheta, z, u^{(k)})) \prod_{i \neq k} Q_{\theta, \vartheta, z}(du^{(i)}) \frac{r_{u^{(k)}}(\theta, \vartheta, z)}{\sum_{i=1}^N r_{u^{(i)}}(\theta, \vartheta, z)} \\ \pi^\varphi(d\xi) &= \hat{\pi}_0^{\varphi_0}(d(\theta, \vartheta, z, u^{(k)})) \prod_{i \neq k} Q_{\theta, \vartheta, z}(du^{(i)}) \frac{1}{N}.\end{aligned}$$

Therefore, for  $c = 1$ , we have

$$\begin{aligned}\hat{r}(\xi) &= \frac{\pi^\varphi(d\xi)}{\pi(d\xi)} = \frac{\hat{\pi}_0^{\varphi_0}(d(\theta, \vartheta, z, u^{(k)})) \prod_{i \neq k} Q_{\theta, \vartheta, z}(du^{(i)})}{\hat{\pi}_0(d(\theta, \vartheta, z, u^{(k)})) \prod_{i \neq k} Q_{\theta, \vartheta, z}(du^{(i)})} \frac{1}{\frac{r_{u^{(k)}}(\theta, \vartheta, z)}{\frac{1}{N} \sum_{i=1}^N r_{u^{(i)}}(\theta, \vartheta, z)}} \\ &= r_u^N(\theta, \vartheta, z),\end{aligned}$$

When  $c = 2$ , we use the relation in (7) to obtain

$$\begin{aligned}\hat{r}(\theta, \vartheta, z, u, k, 2) &= \hat{r}(\vartheta, \theta, z', u', k, 1)^{-1} \\ &= [r_{u'}^N(\vartheta, \theta, z')]^{-1}.\end{aligned}$$

□

## A.2 Proof of Theorem 2

*Proof of Theorem 2.* We start by noticing that the expression for the Dirichlet form associated with  $\hat{P}^N$  can be rewritten in either of the following simplified forms

$$\begin{aligned}\mathcal{E}_{\hat{P}^N}(f) &= \frac{1}{2} \int \pi(d\theta) \int_{\mathfrak{U} \times \llbracket N \rrbracket} Q_1^N(\theta, d(\vartheta, u, k)) \min\{1, r_u^N(x, y)\} [f(\theta) - f(\vartheta)]^2 \\ &= \frac{1}{2} \int \pi(d\theta) \int_{\mathfrak{U} \times \llbracket N \rrbracket} Q_2^N(\theta, d(\vartheta, u, k)) \min\{1, 1/r_u^N(\vartheta, \theta)\} [f(\theta) - f(\vartheta)]^2.\end{aligned}$$

The expression on the first line turns out to be particularly convenient. A well known result from the convex order literature states that for any  $n \geq 2$  exchangeable random variables  $Z_1, \dots, Z_n$  and any convex function  $f$  we have  $\mathbb{E}[f(n^{-1} \sum_{i=1}^n Z_i)] \leq \mathbb{E}[f((n-1)^{-1} \sum_{i=1}^{n-1} Z_i)]$  whenever the expectations exist [Müller and Stoyan, 2002, Corollary 1.5.24]. The two sums are said to be convex ordered. Now since  $a \mapsto -\min\{1, a\}$  is convex we deduce that for any  $N \geq 1$ ,  $\theta, \vartheta \in \Theta$ ,

$$\int_{\mathfrak{U}^N} Q_{\theta, \vartheta}^N(du) \min\{1, r_u^N(\theta, \vartheta)\} \leq \int_{\mathfrak{U}^{N+1}} Q_{\theta, \vartheta}^{N+1}(du) \min\{1, r_u^{N+1}(\theta, \vartheta)\} \quad (39)$$

where  $Q_{\theta, \vartheta}^N(du) := \prod_{i=1}^N Q_{\theta, \vartheta}(du^{(i)})$ , and consequently for any  $f \in L^2(\Theta, \pi)$  and  $N \geq 1$

$$\mathcal{E}_{\hat{P}^{N+1}}(f) \leq \mathcal{E}_{\hat{P}^N}(f).$$

All the monotonicity properties follow from Tierney [1998] since  $\hat{P}^N$  and  $\hat{P}^{N+1}$  are  $\pi$ -reversible. The comparisons to  $P$  follow from the application of Jensen's inequality to  $a \mapsto \min\{1, a\}$ , which leads for any  $\theta, \vartheta \in \Theta$  to

$$\int_{\mathfrak{U}} Q_{\theta, \vartheta}^N(du) \min\{1, r_u^N(\theta, \vartheta)\} \leq \min\{1, r(\theta, \vartheta)\},$$

and again using the results of Tierney [1998].

□

## B Proofs for Section 3

We first prove Theorem 3 which establishes the expression for the acceptance ratio of MHAAR-RB for the multiple latent variable model. Then, we prove the correctness of the delayed rejection algorithm given in Section 3.1.

### B.1 Acceptance ratio of Algorithm 3

For the multiple latent variable model in Section 3.1, recall the joint density

$$\pi(\theta, z) \propto \eta(\theta) \prod_{t=1}^T \gamma_{t,\theta}(z_t). \quad (40)$$

Define  $C_\theta = \prod_{t=1}^T \int_{\mathbf{Z}} \gamma_{t,\theta}(z) dz$ , and  $C = \int_{\Theta} \eta(\theta) C_\theta d\theta$  so that the marginal density is  $\pi(\theta) = \eta(\theta) C_\theta / C$  and the conditional density of the latent variables is

$$\pi_\theta(z) := \frac{\pi(\theta, z)}{\pi(\theta)} = \frac{\prod_{t=1}^T \gamma_{t,\theta}(z_t)}{C_\theta}. \quad (41)$$

Furthermore, for any  $\theta, \vartheta \in \Theta^2$ ,  $t \geq 1$ , and  $z, v \in \mathbf{Z}^2$ , define

$$w_{t,\theta,\vartheta}(v) := \frac{\gamma_{t,\theta,\vartheta}(v)}{q_{t,\theta,\vartheta}(v)}, \quad \text{and} \quad \lambda_{t,\theta,\vartheta}(z, u) := \frac{\gamma_{t,\theta,\vartheta}(z)}{\gamma_{t,\theta}(z)} \frac{\gamma_{t,\vartheta}(u)}{\gamma_{t,\theta,\vartheta}(u)}.$$

We need the following preparatory lemmas for the proofs in this section.

**Lemma 1.** *For any  $\theta, \vartheta \in \Theta$ ,  $\mathbf{v} \in \mathbf{Z}^{MT}$ ,  $\mathbf{k} \in \llbracket M \rrbracket^T$ , and , we have*

$$\frac{(\pi_\theta \otimes \Phi_{\theta,\vartheta})^{\mathbf{s}_{1,\mathbf{k}}}(\mathbf{d}\mathbf{v}) b_{\theta,\vartheta}(\mathbf{k} | \mathbf{s}_{1,\mathbf{k}}(\mathbf{v}))}{(\pi_\theta \otimes \Phi_{\theta,\vartheta})(\mathbf{d}\mathbf{v}) b_{\theta,\vartheta}(\mathbf{k} | \mathbf{v})} = \prod_{t=1}^T \frac{\gamma_{t,\theta}(v_t^{(k_t)})}{\gamma_{t,\theta,\vartheta}(v_t^{(k_t)})} \frac{\gamma_{t,\vartheta}(v_t^{(1)})}{\gamma_{t,\theta}(v_t^{(1)})},$$

where  $\mathbf{s}_{1,\mathbf{k}}$ ,  $\Phi_{\theta,\vartheta}$ ,  $b_{\theta,\vartheta}$  are defined in (19), (17), and in (25), respectively.

*Proof of Lemma 1.* The denominator is equal to

$$\begin{aligned} (\pi_\theta \otimes \Phi_{\theta,\vartheta})(\mathbf{d}\mathbf{v}) b_{\theta,\vartheta}(\mathbf{k} | \mathbf{v}) &= \pi_\theta(v^{(1)}) \left[ \prod_{t=1}^T \prod_{i=2}^M q_{t,\theta,\vartheta}(v_t^{(i)}) \right] \prod_{t=1}^T \frac{w_{t,\theta,\vartheta}(v_t^{(k_t)})}{\sum_{i=1}^M w_{t,\theta,\vartheta}(v_t^{(i)})} \\ &= \frac{1}{C_\theta} \left[ \prod_{t=1}^T \frac{\gamma_{t,\theta}(v_t^{(1)})}{q_{t,\theta,\vartheta}(v_t^{(1)})} \prod_{i=1}^M q_{t,\theta,\vartheta}(v_t^{(i)}) \right] \prod_{t=1}^T \frac{w_{t,\theta,\vartheta}(v_t^{(k_t)})}{\sum_{i=1}^M w_{t,\theta,\vartheta}(v_t^{(i)})} \\ &= \frac{1}{C_\theta} \prod_{t=1}^T \frac{\gamma_{t,\theta}(v_t^{(1)}) \gamma_{t,\theta}(v_t^{(k_t)})}{q_{t,\theta,\vartheta}(v_t^{(1)}) q_{t,\theta,\vartheta}(v_t^{(k_t)})} \frac{\prod_{i=1}^M q_{t,\theta,\vartheta}(v_t^{(i)})}{\sum_{i=1}^M w_{t,\theta,\vartheta}(v_t^{(i)})} \prod_{t=1}^T \frac{\gamma_{t,\vartheta}(v_t^{(k_t)})}{\gamma_{t,\theta}(v_t^{(k_t)})}. \end{aligned}$$

The numerator is obtained by swapping  $v^{(1)}$  and  $v^{(\mathbf{k})}$  in the denominator, therefore it is equal to

$$(\pi_\theta \otimes \Phi_{\theta,\vartheta})^{\mathbf{s}_{1,\mathbf{k}}}(\mathbf{d}\mathbf{v}) b_{\theta,\vartheta}(\mathbf{k} | \mathbf{s}_{1,\mathbf{k}}(\mathbf{v})) = \frac{1}{C_\theta} \prod_{t=1}^T \frac{\gamma_{t,\theta}(v_t^{(1)}) \gamma_{t,\theta}(v_t^{(k_t)})}{q_{t,\theta,\vartheta}(v_t^{(1)}) q_{t,\theta,\vartheta}(v_t^{(k_t)})} \frac{\prod_{i=1}^M q_{t,\theta,\vartheta}(v_t^{(i)})}{\sum_{i=1}^M w_{t,\theta,\vartheta}(v_t^{(i)})} \prod_{t=1}^T \frac{\gamma_{t,\vartheta}(v_t^{(1)})}{\gamma_{t,\theta}(v_t^{(1)})}$$

Taking the ratio yields the result.  $\square$

**Lemma 2.** For any  $\theta, \vartheta \in \Theta$ ,  $\mathbf{v} \in \mathbb{Z}^{MT}$ , and  $\mathbf{k} \in \llbracket M \rrbracket^T$ , we have

$$b_{\theta, \vartheta}^{(1)}(\mathbf{k}|\mathbf{v}) = b_{\theta, \vartheta}(\mathbf{k}|\mathbf{v}) \frac{r_{v^{(1)}, v^{(\mathbf{k})}}(\theta, \vartheta)}{r_{\mathbf{1}, \mathbf{v}}(\theta, \vartheta)}$$

where  $b_{\theta, \vartheta}$  and  $b_{\theta, \vartheta}^{(1)}$  are defined in (18) and (25), and  $r_{v^{(1)}, v^{(\mathbf{k})}}(\theta, \vartheta)$  and  $r_{\mathbf{1}, \mathbf{v}}(\theta, \vartheta)$  are defined in (20) and (21), respectively.

*Proof of Lemma 2.* We prove simply by showing that the ratio

$$\begin{aligned} \frac{b_{\theta, \vartheta}^{(1)}(\mathbf{k}|\mathbf{v})}{b_{\theta, \vartheta}(\mathbf{k}|\mathbf{v})} &= \prod_{t=1}^T \frac{\frac{w_{t, \theta, \vartheta}(v_t^{(k_t)}) \lambda_{t, \theta, \vartheta}(v_t^{(1)}, v_t^{(k_t)})}{\sum_{j=1}^M w_{t, \theta, \vartheta}(v_t^{(j)}) \lambda_{t, \theta, \vartheta}(v_t^{(1)}, v_t^{(j)})}}{\frac{w_{t, \theta, \vartheta}(v_t^{(k_t)})}{\sum_{j=1}^M w_{t, \theta, \vartheta}(v_t^{(j)})}} \\ &= \prod_{t=1}^T \frac{\sum_{j=1}^M w_{t, \theta, \vartheta}(v_t^{(j)})}{\sum_{j=1}^M w_{t, \theta, \vartheta}(v_t^{(j)}) \lambda_{t, \theta, \vartheta}(v_t^{(1)}, v_t^{(j)})} \prod_{t=1}^T \frac{w_{t, \theta, \vartheta}(v_t^{(k_t)}) \lambda_{t, \theta, \vartheta}(v_t^{(1)}, v_t^{(k_t)})}{w_{t, \theta, \vartheta}(v_t^{(k_t)})} \\ &= \frac{\eta(\vartheta) q(\vartheta, \theta)}{\eta(\theta) q(\theta, \vartheta)} \prod_{t=1}^T \lambda_{t, \theta, \vartheta}(v_t^{(1)}, v_t^{(k_t)}) r_{\mathbf{1}, \mathbf{v}}(\theta, \vartheta)^{-1} \\ &= \frac{r_{v^{(1)}, v^{(\mathbf{k})}}(\theta, \vartheta)}{r_{\mathbf{1}, \mathbf{v}}(\theta, \vartheta)} \end{aligned}$$

as claimed.  $\square$

The next lemma can be verified by inspection and therefore will be left without a proof.

**Lemma 3.** For any  $\theta, \vartheta \in \Theta$ ,  $\mathbf{v} \in \mathbb{Z}^{MT}$ , and  $\mathbf{k} \in \llbracket M \rrbracket^T$ , we have  $r_{\mathbf{k}, \mathbf{v}}(\theta, \vartheta) = r_{\mathbf{1}, \mathbf{s}_{\mathbf{1}, \mathbf{k}}(\mathbf{v})}(\theta, \vartheta)$ .

Now we prove Theorem 3 using the lemmas above.

*Proof Theorem 3.* Recalling the notation in Section 3.1, the joint distribution  $\dot{\pi}$  for  $\xi = (\theta, \vartheta, \mathbf{v}, \mathbf{k}, c)$  can be written in compact form as

$$\dot{\pi}(d\xi) = \frac{1}{2} \pi(d\theta) q(\theta, d\vartheta) \left[ \mathbb{I}_1(c) (\pi_\theta \otimes \Phi_{\theta, \vartheta})(d\mathbf{v}) b_{\theta, \vartheta}^{(1)}(\mathbf{k}|\mathbf{v}) + \mathbb{I}_2(c) (\pi_\theta \otimes \Phi_{\vartheta, \theta})(d\mathbf{v}) b_{\vartheta, \theta}(\mathbf{k}|\mathbf{v}) \right],$$

and the proposed involution is  $\varphi(\theta, \vartheta, \mathbf{v}, \mathbf{k}, c) = (\vartheta, \theta, \mathbf{s}_{\mathbf{1}, \mathbf{k}}(\mathbf{v}), \mathbf{k}, 3 - c)$ . When  $c = 1$ ,  $\dot{\pi}(\xi)$  is

$$\begin{aligned} \frac{\dot{\pi}^\varphi(d\xi)}{\dot{\pi}(d\xi)} &= \frac{q(\vartheta, \theta) \pi(\vartheta)}{q(\theta, \vartheta) \pi(\theta)} \frac{(\pi_\theta \otimes \Phi_{\theta, \vartheta})^{\mathbf{s}_{\mathbf{1}, \mathbf{k}}}(\mathbf{v}) b_{\theta, \vartheta}(\mathbf{k}|\mathbf{s}_{\mathbf{1}, \mathbf{k}}(\mathbf{v}))}{(\pi_\theta \otimes \Phi_{\theta, \vartheta})(d\mathbf{v}) b_{\theta, \vartheta}^{(1)}(\mathbf{k}|\mathbf{v})} \\ &= \frac{q(\vartheta, \theta) \pi(\vartheta)}{q(\theta, \vartheta) \pi(\theta)} \frac{\pi_\vartheta(dv^{(\mathbf{k})})}{\pi_\theta(dv^{(\mathbf{k})})} \frac{(\pi_\theta \otimes \Phi_{\theta, \vartheta})^{\mathbf{s}_{\mathbf{1}, \mathbf{k}}}(\mathbf{v}) b_{\theta, \vartheta}(\mathbf{k}|\mathbf{s}_{\mathbf{1}, \mathbf{k}}(\mathbf{v}))}{(\pi_\theta \otimes \Phi_{\theta, \vartheta})(d\mathbf{v}) b_{\theta, \vartheta}(\mathbf{k}|\mathbf{v})} \frac{r_{\mathbf{1}, \mathbf{v}}(\theta, \vartheta)}{r_{v^{(1)}, v^{(\mathbf{k})}}(\theta, \vartheta)} \\ &= \frac{q(\vartheta, \theta) \pi(\vartheta)}{q(\theta, \vartheta) \pi(\theta)} \frac{\pi_\vartheta(v^{(\mathbf{k})})}{\pi_\theta(v^{(\mathbf{k})})} \prod_{t=1}^T \frac{\gamma_{t, \theta}(v_t^{(k_t)})}{\gamma_{t, \theta, \vartheta}(v_t^{(k_t)})} \frac{\gamma_{t, \theta, \vartheta}(v_t^{(1)})}{\gamma_{t, \theta}(v_t^{(1)})} \frac{1}{r_{v^{(1)}, v^{(\mathbf{k})}}(\theta, \vartheta)} r_{\mathbf{1}, \mathbf{v}}(\theta, \vartheta) \\ &= \frac{q(\vartheta, \theta) \eta(\vartheta)}{q(\theta, \vartheta) \eta(\theta)} \prod_{t=1}^T \frac{\gamma_{t, \vartheta}(v_t^{(k_t)})}{\gamma_{t, \theta}(v_t^{(k_t)})} \frac{\gamma_{t, \theta}(v_t^{(k_t)})}{\gamma_{t, \theta, \vartheta}(v_t^{(k_t)})} \frac{\gamma_{t, \theta, \vartheta}(v_t^{(1)})}{\gamma_{t, \theta}(v_t^{(1)})} \frac{1}{r_{v^{(1)}, v^{(\mathbf{k})}}(\theta, \vartheta)} r_{\mathbf{1}, \mathbf{v}}(\theta, \vartheta) \\ &= r_{\mathbf{1}, \mathbf{v}}(\theta, \vartheta), \end{aligned}$$



where we have used Lemma 2 for the second line, Lemma 1 for the third line, (40) for the fourth line, and the definition of acceptance ratio  $r_{v^{(1)}, v^{(k)}}(\theta, \vartheta)$  in (20) for the last line.

For  $c = 2$ , upon using (7), we write

$$\begin{aligned} \mathring{r}(\theta, \vartheta, \mathbf{v}, \mathbf{k}, 2) &= \mathring{r}(\vartheta, \theta, \mathfrak{s}_{\mathbf{1}, \mathbf{k}}(\mathbf{v}), \mathbf{k}, 1)^{-1} \\ &= r_{\mathbf{1}, \mathfrak{s}_{\mathbf{1}, \mathbf{k}}(\mathbf{v})}(\vartheta, \theta)^{-1} \\ &= r_{\mathbf{k}, \mathbf{v}}(\vartheta, \theta)^{-1}, \end{aligned}$$

where the last line is due to  $\mathfrak{s}_{\mathbf{1}, \mathbf{k}}(\mathbf{v})^{(1)} = v^{(k)}$  and Lemma 3.  $\square$

## B.2 Delayed rejection step for Algorithm 3

When the delayed rejection step is included in Algorithm 3, the algorithm targets the modified joint distribution for  $\check{\xi} := (\xi, \mathbf{l}, \mathbf{l}')$  defined as

$$\tilde{\pi}(\mathrm{d}\check{\xi}) = \mathring{\pi}(\mathrm{d}\xi) [\mathbb{I}_1(c) b^{\text{ref}}(\mathbf{l}, \mathbf{l}' | \xi) + \mathbb{I}_2(c) b^{\text{ref}}(\mathbf{l}', \mathbf{l} | \varphi(\xi))]$$

where  $\xi = (\theta, \vartheta, \mathbf{v}, \mathbf{k}, c)$  is as in Section B.1, and, conditional on  $\xi$ , the joint probability distribution of  $\mathbf{l}, \mathbf{l}' \in \llbracket M \rrbracket^T$  is given by

$$b^{\text{ref}}(\mathbf{l}, \mathbf{l}' | \xi) := b_{\theta, \vartheta}^{\text{ref}, (1)}(\mathbf{l} | \mathbf{v}) b_{\vartheta, \theta}^{\text{ref}, (2)}(\mathbf{l}' | \mathfrak{s}_{\mathbf{1}, \mathbf{k}}(\mathbf{v}))$$

with the individual probabilities defined as

$$b_{\theta, \vartheta}^{\text{ref}, (1)}(\mathbf{l} | \mathbf{v}) := \prod_{t=1}^T \frac{\gamma_{t, \theta}(v_t^{(l_t)}) / q_{t, \theta, \vartheta}(v_t^{(l_t)})}{\sum_{i=1}^N \gamma_{t, \theta}(v_t^{(i)}) / q_{t, \theta, \vartheta}(v_t^{(i)})}, \quad b_{\vartheta, \theta}^{\text{ref}, (2)}(\mathbf{l}' | \mathbf{v}) := \prod_{t=1}^T \frac{\gamma_{t, \vartheta}(v_t^{(l'_t)}) / q_{t, \vartheta, \theta}(v_t^{(l'_t)})}{\sum_{i=1}^N \gamma_{t, \vartheta}(v_t^{(i)}) / q_{t, \vartheta, \theta}(v_t^{(i)})}.$$

These are simply the selection probabilities of  $\gamma_{t, \theta}$ -invariant cSMC kernels when the proposed values are sampled from  $q_{t, \theta, \vartheta}$  or  $q_{t, \vartheta, \theta}$ , respectively.

The algorithm can be thought of as a two-stage delayed rejection algorithm, where the stage one move corresponds to the regular MHAAR update and stage two move is executed only if the move in stage one is rejected. We note that, in practice, the pair  $\mathbf{l}, \mathbf{l}'$  do not play any role in the implementation of the first stage. Moreover, in the implementation of the second stage, one only needs to sample  $\mathbf{l}$  to propose  $v^{(1)}$  for the latent variable;  $\mathbf{l}'$  is, again, not needed.

The mentioned two stages of the delayed rejection algorithm are given below.

1. In the first stage, MHAAR attempts a transition for the joint variable  $\check{\xi} = (\xi, \mathbf{l}, \mathbf{l}')$  as

$$\check{\varphi}_1(\xi, \mathbf{l}, \mathbf{l}') := (\varphi(\xi), \mathbf{l}', \mathbf{l}).$$

As  $\check{\varphi}_1$  is an involution, it yields the acceptance ratio

$$\begin{aligned} \check{r}_1(\check{\xi}) &:= \frac{\tilde{\pi}^{\check{\varphi}_1}(\mathrm{d}\check{\xi})}{\tilde{\pi}(\mathrm{d}\check{\xi})} \\ &= \frac{\mathring{\pi}^{\varphi}(\mathrm{d}\xi)}{\mathring{\pi}(\mathrm{d}\xi)} \frac{b^{\text{ref}}(\mathbf{l}, \mathbf{l}' | \varphi \circ \varphi(\xi))}{b^{\text{ref}}(\mathbf{l}, \mathbf{l}' | \xi)} \\ &= \frac{\mathring{\pi}^{\varphi}(\mathrm{d}\xi)}{\mathring{\pi}(\mathrm{d}\xi)} = \mathring{r}(\xi), \end{aligned} \tag{42}$$



which is exactly the same acceptance ratio as we would have without the delayed rejection step. Note that, neither the acceptance ratio nor the variables carried on to the next iteration depend on the additional variables  $\mathbf{l}$  or  $\mathbf{l}'$ . Therefore, the variables  $\mathbf{l}$  and  $\mathbf{l}'$  need not be sampled prior to the delayed rejection step.

2. The second stage corresponds to proposing a transformation of  $\check{\xi} = (\xi, \mathbf{l}, \mathbf{l}')$ , recalling that  $\xi = (\theta, \vartheta, \mathbf{v}, \mathbf{k}, c)$ , with the following involution:

$$\check{\varphi}_2(\xi, \mathbf{l}, \mathbf{l}') := (\theta, \vartheta, \mathbf{s}_{\mathbf{l}, \mathbf{l}'}(\mathbf{v}), \mathbf{r}_1(\mathbf{k}), c, \mathbf{l}, \mathbf{r}_1(\mathbf{l}'))$$

where, for any  $\mathbf{k}, \mathbf{l} \in \llbracket M \rrbracket^T$ , we define  $\mathbf{r}_1(\mathbf{k}) : \llbracket M \rrbracket^T \mapsto \llbracket M \rrbracket^T$  as,

$$[\mathbf{r}_1(\mathbf{k})]_i = \begin{cases} l_i, & k_i = 1 \\ 1, & k_i = l_i \\ k_i, & \text{otherwise} \end{cases}, \quad i = 1, \dots, T, \quad (43)$$

That is,  $\mathbf{r}_1(\mathbf{k})$  is the set of indices of the elements of  $v^{(\mathbf{k})}$  once  $v^{(1)}$  and  $v^{(\mathbf{l})}$  have been swapped in  $\mathbf{v}$ , so that  $[\mathbf{s}_{\mathbf{l}, \mathbf{l}'}(\mathbf{v})]^{(\mathbf{k})} = v^{(\mathbf{r}_1(\mathbf{k}))}$ . We note that, the operator  $\mathbf{r}_1$  is merely introduced to establish the correctness of the algorithm and in practice does not need to be implemented. Crucially for our analysis, it can be checked that, for any  $\mathbf{l} \in \llbracket M \rrbracket^T$ , the operator  $\mathbf{r}_1(\cdot)$  is an involution, resulting in  $\check{\varphi}_2$  also being an involution. This enables us to cast the delayed rejection scheme in our framework. the acceptance ratio in the second stage can be written as

$$\check{r}_2(\check{\xi}) = \frac{\tilde{\pi}^{\check{\varphi}_2}(\mathrm{d}\check{\xi})}{\tilde{\pi}(\mathrm{d}\check{\xi})} \frac{1 - \min\{1, \check{r}_1 \circ \check{\varphi}_2(\check{\xi})\}}{1 - \min\{1, \check{r}_1(\check{\xi})\}} \quad (44)$$

**Theorem 6.** *The acceptance ratio in (44) is equal to*

$$\check{r}_2(\check{\xi}) = \begin{cases} \frac{1 - \min\{0, r_{\mathbf{l}, \mathbf{v}}(\theta, \vartheta)\}}{1 - \min\{0, r_{\mathbf{l}, \mathbf{v}}(\theta, \vartheta)\}}, & c = 1; \\ \frac{1 - \min\{0, 1/r_{\mathbf{k}, \mathbf{v}}(\vartheta, \theta)\}}{1 - \min\{0, 1/r_{\mathbf{l}, \mathbf{v}}(\vartheta, \theta)\}}, & c = 2. \end{cases}$$

Moreover, when  $\gamma_{t, \theta, \vartheta} = \gamma_{t, \theta}$  for all  $t, \theta, \vartheta$ , the acceptance ratio simplifies to  $\check{r}_2(\check{\xi}) = 1$ .

In the proof of Theorem 6, we will make use of the following lemmas.

**Lemma 4.** *For any  $\theta, \vartheta \in \Theta$ ,  $\mathbf{v} \in Z^{MT}$ , and  $\mathbf{k}, \mathbf{l} \in \llbracket M \rrbracket^T$ , we have the following facts.*

- $[\mathbf{s}_{\mathbf{l}, \mathbf{l}'}(\mathbf{v})]^{(\mathbf{r}_1(\mathbf{k}))} = v^{(\mathbf{k})}$ .
- The following equalities hold

$$(\pi_\theta \otimes \Phi_{\theta, \vartheta})(\mathrm{d}\mathbf{v}) b_{\theta, \vartheta}^{\text{ref}, (1)}(\mathbf{k}|\mathbf{v}) = (\pi_\theta \otimes \Phi_{\theta, \vartheta})^{\mathbf{s}_{\mathbf{l}, \mathbf{k}}}(\mathrm{d}\mathbf{v}) b_{\theta, \vartheta}^{\text{ref}, (1)}(\mathbf{k}|\mathbf{s}_{\mathbf{l}, \mathbf{k}}(\mathbf{v})) \quad (45)$$

$$(\pi_\theta \otimes \Phi_{\vartheta, \theta})(\mathrm{d}\mathbf{v}) b_{\vartheta, \theta}^{\text{ref}, (2)}(\mathbf{k}|\mathbf{v}) = (\pi_\theta \otimes \Phi_{\vartheta, \theta})^{\mathbf{s}_{\mathbf{l}, \mathbf{k}}}(\mathrm{d}\mathbf{v}) b_{\vartheta, \theta}^{\text{ref}, (2)}(\mathbf{k}|\mathbf{s}_{\mathbf{l}, \mathbf{k}}(\mathbf{v})) \quad (46)$$

*Proof Lemma 4.* The the first part of Lemma 4 can be verified by inspection. For (45), we write

$$\begin{aligned}
(\pi_\theta \otimes \Phi_{\theta, \vartheta})(d\mathbf{v}) b_{\theta, \vartheta}^{\text{ref}, (1)}(\mathbf{k}|\mathbf{v}) &= \pi_\theta(v^{(1)}) \left[ \prod_{t=1}^T \prod_{i=2}^M q_{t, \theta, \vartheta}(v_t^{(i)}) \right] \prod_{t=1}^T \frac{\gamma_{t, \theta}(v_t^{(k_t)})/q_{t, \theta, \vartheta}(v_t^{(k_t)})}{\sum_{i=1}^N \gamma_{t, \theta}(v_t^{(i)})/q_{t, \theta, \vartheta}(v_t^{(i)})} \\
&= \pi_\theta(v^{(\mathbf{k})}) \left[ \prod_{t=1}^T \frac{\gamma_{t, \theta}(v_t^{(1)})}{\gamma_{t, \theta}(v_t^{(k_t)})} \prod_{i=2}^M q_{t, \theta, \vartheta}(v_t^{(i)}) \right] \prod_{t=1}^T \frac{\gamma_{t, \theta}(v_t^{(k_t)})/q_{t, \theta, \vartheta}(v_t^{(k_t)})}{\sum_{i=1}^N \gamma_{t, \theta}(v_t^{(i)})/q_{t, \theta, \vartheta}(v_t^{(i)})} \\
&= \pi_\theta(v^{(\mathbf{k})}) \left[ \prod_{t=1}^T \frac{\gamma_{t, \theta}(v_t^{(1)})}{\gamma_{t, \theta}(v_t^{(k_t)}) q_{t, \theta, \vartheta}(v_t^{(1)})} \prod_{i=1}^M q_{t, \theta, \vartheta}(v_t^{(i)}) \right] \prod_{t=1}^T \frac{\gamma_{t, \theta}(v_t^{(k_t)})/q_{t, \theta, \vartheta}(v_t^{(k_t)})}{\sum_{i=1}^N \gamma_{t, \theta}(v_t^{(i)})/q_{t, \theta, \vartheta}(v_t^{(i)})} \\
&= \pi_\theta(v^{(\mathbf{k})}) \left[ \prod_{t=1}^T \frac{\gamma_{t, \theta}(v_t^{(k_t)})}{\gamma_{t, \theta}(v_t^{(k_t)}) q_{t, \theta, \vartheta}(v_t^{(k_t)})} \prod_{i=1}^M q_{t, \theta, \vartheta}(v_t^{(i)}) \right] \prod_{t=1}^T \frac{\gamma_{t, \theta}(v_t^{(1)})/q_{t, \theta, \vartheta}(v_t^{(1)})}{\sum_{i=1}^N \gamma_{t, \theta}(v_t^{(i)})/q_{t, \theta, \vartheta}(v_t^{(i)})} \\
&= \pi_\theta(v^{(\mathbf{k})}) \left[ \prod_{t=1}^T \prod_{i \neq k_t}^M q_{t, \theta, \vartheta}(v_t^{(i)}) \right] \prod_{t=1}^T \frac{\gamma_{t, \theta}(v_t^{(1)})/q_{t, \theta, \vartheta}(v_t^{(1)})}{\sum_{i=1}^N \gamma_{t, \theta}(v_t^{(i)})/q_{t, \theta, \vartheta}(v_t^{(i)})} \\
&= (\pi_\theta \otimes \Phi_{\theta, \vartheta})^{\mathbf{s}_{1, \mathbf{k}}}(\mathbf{d}\mathbf{v}) b_{\theta, \vartheta}^{\text{ref}, (1)}(\mathbf{1}|\mathbf{v}) \\
&= (\pi_\theta \otimes \Phi_{\theta, \vartheta})^{\mathbf{s}_{1, \mathbf{k}}}(\mathbf{d}\mathbf{v}) b_{\theta, \vartheta}^{\text{ref}, (1)}(\mathbf{k}|\mathbf{s}_{1, \mathbf{k}}(\mathbf{v})),
\end{aligned}$$

hence (45) is shown. We prove (46) using the same steps above, replacing  $\Phi_{\theta, \vartheta}(\cdot)$ ,  $q_{t, \theta, \vartheta}(\cdot)$ , and  $b_{\theta, \vartheta}^{\text{ref}, (1)}(\cdot)$  by  $\Phi_{\vartheta, \theta}(\cdot)$ ,  $q_{t, \vartheta, \theta}(\cdot)$ , and  $b_{\vartheta, \theta}^{\text{ref}, (2)}$ , respectively.  $\square$

The following lemma can be verified by inspection, hence we skip a formal proof.

**Lemma 5.** Suppose  $\gamma_{t, \theta, \vartheta} = \gamma_{t, \theta}$ , for all  $t \geq 1$ , and  $\theta, \vartheta \in \Theta$ . Then, for any  $\mathbf{v} \in \mathbb{Z}^{MT}$ , and  $\mathbf{k}, \mathbf{l} \in \llbracket M \rrbracket^T$ , we have  $r_{\mathbf{k}, \mathbf{v}}(\theta, \vartheta) = r_{\mathbf{l}, \mathbf{v}}(\theta, \vartheta)$ .

We proceed to the proof of Theorem 6.

*Proof of Theorem 6.* First, we show that the first ratio in (44) is equal to 1. Indeed, for  $c = 1$ ,

$$\begin{aligned}
\frac{\tilde{\pi}^{\varphi_2}(\tilde{d}\tilde{\xi})}{\tilde{\pi}(\tilde{d}\tilde{\xi})} &= r(\theta, \vartheta) \frac{(\pi_\theta \otimes \Phi_{\theta, \vartheta})^{\mathbf{s}_{1, \mathbf{l}}}(\mathbf{d}\mathbf{v}) b_{\theta, \vartheta}^{(1)}(\mathbf{r}_1(\mathbf{k})|\mathbf{s}_{1, \mathbf{l}}(\mathbf{v})) b_{\theta, \vartheta}^{\text{ref}, (1)}(\mathbf{l}|\mathbf{s}_{1, \mathbf{l}}(\mathbf{v})) b_{\vartheta, \theta}^{\text{ref}, (2)}(\mathbf{r}_1(\mathbf{l}')|\mathbf{s}_{1, \mathbf{l}}(\mathbf{v}))}{(\pi_\theta \otimes \Phi_{\theta, \vartheta})(\mathbf{d}\mathbf{v}) b_{\theta, \vartheta}^{(1)}(\mathbf{k}|\mathbf{v}) b_{\theta, \vartheta}^{\text{ref}, (1)}(\mathbf{l}|\mathbf{v}) b_{\vartheta, \theta}^{\text{ref}, (2)}(\mathbf{l}'|\mathbf{v})} \\
&= r(\theta, \vartheta) \frac{(\pi_\theta \otimes \Phi_{\theta, \vartheta})^{\mathbf{s}_{1, \mathbf{l}}}(\mathbf{d}\mathbf{v}) b_{\theta, \vartheta}^{\text{ref}, (1)}(\mathbf{l}|\mathbf{s}_{1, \mathbf{l}}(\mathbf{v})) b_{\theta, \vartheta}^{(1)}(\mathbf{r}_1(\mathbf{k})|\mathbf{s}_{1, \mathbf{l}}(\mathbf{v})) b_{\vartheta, \theta}^{\text{ref}, (2)}(\mathbf{r}_1(\mathbf{l}')|\mathbf{s}_{1, \mathbf{l}}(\mathbf{v}))}{(\pi_\theta \otimes \Phi_{\theta, \vartheta})(\mathbf{d}\mathbf{v}) b_{\theta, \vartheta}^{\text{ref}, (1)}(\mathbf{l}|\mathbf{v}) b_{\theta, \vartheta}^{(1)}(\mathbf{k}|\mathbf{v}) b_{\vartheta, \theta}^{\text{ref}, (2)}(\mathbf{l}'|\mathbf{v})},
\end{aligned}$$

and all of the ratios are equal to 1, due to Lemma 4. For  $c = 2$ ,

$$\begin{aligned}
\frac{\tilde{\pi}^{\varphi_2}(\tilde{d}\tilde{\xi})}{\tilde{\pi}(\tilde{d}\tilde{\xi})} &= r(\theta, \vartheta) \frac{(\pi_\theta \otimes \Phi_{\vartheta, \theta})^{\mathbf{s}_{1, \mathbf{l}}}(\mathbf{d}\mathbf{v}) b_{\vartheta, \theta}(\mathbf{r}_1(\mathbf{k})|\mathbf{s}_{1, \mathbf{l}}(\mathbf{v})) b_{\vartheta, \theta}^{\text{ref}, (2)}(\mathbf{l}|\mathbf{s}_{1, \mathbf{l}}(\mathbf{v})) b_{\vartheta, \theta}^{\text{ref}, (1)}(\mathbf{r}_1(\mathbf{l}')|\mathbf{s}_{1, \mathbf{l}}(\mathbf{v}))}{(\pi_\theta \otimes \Phi_{\vartheta, \theta})(\mathbf{d}\mathbf{v}) b_{\vartheta, \theta}(\mathbf{k}|\mathbf{v}) b_{\vartheta, \theta}^{\text{ref}, (2)}(\mathbf{l}|\mathbf{v}) b_{\vartheta, \theta}^{\text{ref}, (1)}(\mathbf{l}'|\mathbf{v})} \\
&= r(\theta, \vartheta) \frac{(\pi_\theta \otimes \Phi_{\vartheta, \theta})^{\mathbf{s}_{1, \mathbf{l}}}(\mathbf{d}\mathbf{v}) b_{\vartheta, \theta}^{\text{ref}, (2)}(\mathbf{l}|\mathbf{s}_{1, \mathbf{l}}(\mathbf{v})) b_{\vartheta, \theta}(\mathbf{r}_1(\mathbf{k})|\mathbf{s}_{1, \mathbf{l}}(\mathbf{v})) b_{\vartheta, \theta}^{\text{ref}, (1)}(\mathbf{r}_1(\mathbf{l}')|\mathbf{s}_{1, \mathbf{l}}(\mathbf{v}))}{(\pi_\theta \otimes \Phi_{\vartheta, \theta})(\mathbf{d}\mathbf{v}) b_{\vartheta, \theta}^{\text{ref}, (2)}(\mathbf{l}|\mathbf{v}) b_{\vartheta, \theta}(\mathbf{k}|\mathbf{v}) b_{\vartheta, \theta}^{\text{ref}, (1)}(\mathbf{l}'|\mathbf{v})},
\end{aligned}$$

and all of the ratios are equal to 1, again, due to Lemma 4.

The second ratio in (44) for any choice of  $\gamma_{t,\theta,\vartheta}$  is equal to as is equal to 1 when  $\gamma_{t,\theta,\vartheta} = \gamma_{t,\theta}$ . we can write the ratio of rejection probabilities ,

$$\begin{aligned} \frac{1 - \min \{1, \check{r}_1 \circ \check{\varphi}_2(\check{\xi})\}}{1 - \min \{1, \check{r}_1(\check{\xi})\}} &= \frac{1 - \min \{1, \check{r}(\theta, \vartheta, \mathbf{s}_{1,1}(\mathbf{v}), \mathbf{r}_1(\mathbf{k}), c)\}}{1 - \min \{1, \check{r}(\theta, \vartheta, \mathbf{v}, \mathbf{k}, c)\}} \\ &= \begin{cases} \frac{1 - \min \{0, r_{1,\mathbf{v}}(\theta, \vartheta)\}}{1 - \min \{0, r_{1,\mathbf{v}}(\theta, \vartheta)\}}, & c = 1; \\ \frac{1 - \min \{0, 1/r_{\mathbf{k},\mathbf{v}}(\vartheta, \theta)\}}{1 - \min \{0, 1/r_{1,\mathbf{v}}(\vartheta, \theta)\}}, & c = 2. \end{cases} \end{aligned}$$

where we use (42) in the first line and the second line is due to Lemma 3. When  $\gamma_{t,\theta,\vartheta} = \gamma_{t,\theta}$ , we use Lemma 5 to conclude that both ratios for  $c = 1$  and  $c = 2$  simplify to 1.  $\square$

## C Auxiliary results and proofs Section 4

First, we lay out some useful results on SMC, cSMC, for the state-space model defined in Section 4.1.

It is standard that the law of a particle filter with  $M$  particles and multinomial resampling for  $\theta \in \Theta$ ,  $\mathbf{v} \in \mathbb{Z}^{MT}$  and ancestral indices  $a \in \llbracket M \rrbracket^{M(T-1)}$  is [Andrieu et al., 2010]

$$\psi_\theta(d(\mathbf{v}, a)) = \prod_{i=1}^M f_\theta(dv_1^{(i)}) \prod_{t=2}^T \left\{ \prod_{i=1}^M \frac{w_{t-1,\theta}(v_{t-1}^{(a_{t-1}^{(i)})})}{\sum_{j=1}^M w_{t-1,\theta}(v_{t-1}^{(j)})} f_\theta(v_{t-1}^{(a_{t-1}^{(i)})}, dv_t^{(i)}) \right\}.$$

What is important for us is that the marginal distribution  $\psi_\theta(d\mathbf{v})$  has a simple form

$$\psi_\theta(d\mathbf{v}) = \prod_{i=1}^M f_\theta(dv_1^{(i)}) \prod_{t=2}^T \left\{ \prod_{i=1}^M \frac{\sum_{j=1}^M w_{t-1,\theta}(v_{t-1}^{(j)}) f_\theta(v_{t-1}^{(j)}, dv_t^{(i)})}{\sum_{j=1}^M w_{t-1,\theta}(v_{t-1}^{(j)})} \right\}.$$

Now, letting  $C_\theta := \ell_\theta(y)$  in (31) (recall  $y = y_{1:T}$ ), and its estimator  $\hat{C}_\theta(\mathbf{v}) := \prod_{t=1}^T \frac{1}{M} \sum_{i=1}^M w_{t,\theta}(v_t^{(i)})$ , we introduce

$$\bar{\psi}_\theta(d\mathbf{v}) := \psi_\theta(d\mathbf{v}) \frac{\hat{C}_\theta(\mathbf{v})}{C_\theta}. \quad (47)$$

We know from Andrieu et al. [2010] that this is a probability distribution, and is a way of justifying that  $\hat{C}_\theta(v)$  is an unbiased estimator of  $C_\theta$ —note that the ancestral history is here integrated out.

The cSMC algorithm is given in Algorithm 6. The joint distribution of  $\mathbf{v} \in \mathbb{Z}^{MT}$  when  $v^{(1)} \sim \pi_\theta(\cdot)$  and  $v^{(\bar{1})}$  is sampled by the cSMC kernel targeting  $\pi_\theta$  can be written as

$$(\pi_\theta \otimes \Phi_\vartheta)(d\mathbf{v}) := \pi_\theta(dv^{(1)}) \prod_{i=2}^M f_\vartheta(dv_1^{(i)}) \prod_{t=2}^T \left\{ \prod_{i=2}^M \frac{\sum_{j=1}^M w_{t-1,\vartheta}(v_{t-1}^{(j)}) f_\vartheta(v_{t-1}^{(j)}, dv_t^{(i)})}{\sum_{j=1}^M w_{t-1,\vartheta}(v_{t-1}^{(j)})} \right\}. \quad (48)$$

Recall the law of the indices used in the backward-sampling procedure in order to draw a path  $v^{(\mathbf{k})}$ ,

$$b_\theta(\mathbf{k}|\mathbf{v}) := \frac{w_{T,\theta}(v_T^{(k_T)})}{\sum_{i=1}^M w_{T,\theta}(v_T^{(i)})} \prod_{t=2}^T \frac{w_{t-1,\theta}(v_{t-1}^{(k_{t-1})}) f_\theta(v_{t-1}^{(k_{t-1})}, v_t^{(k_t)})}{\sum_{i=1}^M w_{t-1,\theta}(v_{t-1}^{(i)}) f_\theta(v_{t-1}^{(i)}, v_t^{(k_t)})}.$$

---

**Algorithm 6:** cSMC( $M, \theta, z$ )

---

**Input:** Number of particles  $M$ , parameter  $\theta$ , current sample  $z$

**Output:** Particles  $\mathbf{v} = v_{1:T}^{(1:M)}$ , new sample  $z'$

```
1 Set  $v_1^{(1)} = z_1$ .
2 for  $i = 2, \dots, M$  do
3   Sample  $v_1^{(i)} \sim f_\theta(\cdot)$ .
4   Compute  $w_1^{(i)} = g_\theta(v_1^{(i)}, y_1)$ .
5 for  $t = 2, \dots, T$  do
6   Set  $v_t^{(1)} = z_t$ .
7   for  $i = 2, \dots, M$  do
8     Sample  $a_{t-1}^{(i)} \sim \mathcal{P}(w_{t-1}^{(1)}, \dots, w_{t-1}^{(M)})$  and  $v_t^{(i)} \sim f_\theta(v_{t-1}^{(a_{t-1}^{(i)})}, \cdot)$ .
9     Compute  $w_t^{(i)} = g_\theta(v_t^{(i)}, y_t)$ .
10 Sample  $k_T \sim \mathcal{P}(w_T^{(1)}, \dots, w_T^{(M)})$  and set  $z'_T = v_T^{(k_T)}$ .
11 for  $t = T-1, \dots, 1$  do
12   for  $i = 1, \dots, M$  do
13     Compute  $\tilde{w}_t^{(i)} = w_t^{(i)} f_\theta(v_t^{(i)}, v_{t+1}^{(k_{t+1})})$ .
14   Sample  $k_t \sim \mathcal{P}(\tilde{w}_t^{(1)}, \dots, \tilde{w}_t^{(M)})$  and set  $z'_t = v_t^{(k_t)}$ .
15 return  $\mathbf{v} = v_{1:T}^{(1:N)}$  and  $z' = z'_{1:T}$ .
```

---

**Lemma 6.** For any  $\theta \in \Theta$  and  $\mathbf{v} \in \mathcal{Z}^{MT}$ ,

$$(\pi_\theta \otimes \Phi_\theta)(d\mathbf{v}) = M^T \bar{\psi}_\theta(d\mathbf{v}) b_\theta(\mathbf{1}|\mathbf{v}).$$

*Proof of Lemma 6.* We check that the ratio

$$\begin{aligned} \frac{(\pi_\theta \otimes \Phi_\theta)(d\mathbf{v})}{\bar{\psi}_\theta(d\mathbf{v})} &= \frac{C_\theta}{\hat{C}_\theta(\mathbf{v})} \frac{\pi_\theta(dv^{(1)}) \prod_{i=2}^M f_\theta(dv_1^{(i)}) \prod_{t=2}^T \left\{ \prod_{i=2}^M \frac{\sum_{j=1}^M w_{t-1, \theta}(v_{t-1}^{(j)}) f_\theta(v_{t-1}^{(j)}, dv_t^{(i)})}{\sum_{j=1}^M w_{t-1, \theta}(v_{t-1}^{(j)})} \right\}}{\prod_{i=1}^M f_\theta(dv_1^{(i)}) \prod_{t=2}^T \left\{ \prod_{i=1}^M \frac{\sum_{j=1}^M w_{t-1, \theta}(v_{t-1}^{(j)}) f_\theta(v_{t-1}^{(j)}, dv_t^{(i)})}{\sum_{j=1}^M w_{t-1, \theta}(v_{t-1}^{(j)})} \right\}} \\ &= \frac{C_\theta}{\prod_{t=1}^T \frac{1}{M} \sum_{i=1}^M w_{\theta, t}(v_t^{(i)})} \frac{\pi_\theta(dv^{(1)})}{f_\theta(dv_1^{(1)})} \frac{\prod_{t=2}^T \sum_{j=1}^M w_{t-1, \theta}(v_{t-1}^{(j)})}{\prod_{t=2}^T \sum_{j=1}^M w_{t-1, \theta}(v_{t-1}^{(j)}) f_\theta(v_{t-1}^{(j)}, dv_t^{(1)})} \\ &= M^T \frac{w_{1, \theta}(v_1^{(1)}) \prod_{t=2}^T f_\theta(v_{t-1}^{(1)}, dv_t^{(1)}) w_{t, \theta}(v_t^{(1)})}{\sum_{i=1}^M w_{1, \theta}(v_1^{(i)}) \prod_{t=2}^T \sum_{j=1}^M w_{t-1, \theta}(v_{t-1}^{(j)}) f_\theta(v_{t-1}^{(j)}, dv_t^{(1)})} \\ &= M^T \frac{w_{1, \theta}(v_1^{(1)})}{\sum_{i=1}^M w_{1, \theta}(v_1^{(i)})} \prod_{t=2}^T \frac{f_\theta(v_{t-1}^{(1)}, dv_t^{(1)}) w_{t, \theta}(v_t^{(1)})}{\sum_{j=1}^M w_{t-1, \theta}(v_{t-1}^{(j)}) f_\theta(v_{t-1}^{(j)}, dv_t^{(1)})} \\ &= M^T b_\theta(\mathbf{1}|\mathbf{v}), \end{aligned}$$

as claimed.  $\square$

The constant  $M^T$  on the right hand side arises from deterministic assignment of indices

$\mathbf{1} = (1, \dots, 1)$  for the conditioned path  $v^{(1)}$  in the cSMC algorithm. Lemmas 7 and 8 can be verified by inspection.

**Lemma 7.** *For any  $\theta \in \Theta$ ,  $\mathbf{v} \in \mathbb{Z}^{MT}$ , and  $\mathbf{k} \in \llbracket M \rrbracket^T$ , and with  $\mathbf{s}_{1,\mathbf{k}}$  defined in (19), we have  $\bar{\psi}_\theta^{\mathbf{s}_{1,\mathbf{k}}}(\mathrm{d}\mathbf{v}) = \bar{\psi}_\theta(\mathrm{d}\mathbf{v})$ .*

**Lemma 8.** *For any  $\theta \in \Theta$ ,  $\mathbf{v} \in \mathbb{Z}^{MT}$ , and  $\mathbf{k} \in \llbracket M \rrbracket^T$ , we have  $b_\theta(\mathbf{k}|\mathbf{s}_{1,\mathbf{k}}(\mathbf{v})) = b_\theta(\mathbf{1}|\mathbf{v})$ .*

Lemmas 6, 7, and 8 lead to the following corollaries which will be useful in the subsequent proofs.

**Corollary 2.** *For any  $\theta \in \Theta$ ,  $\mathbf{v} \in \mathbb{Z}^{MT}$ , and  $\mathbf{k} \in \llbracket M \rrbracket^T$ ,*

$$(\pi_\theta \otimes \Phi_\theta)(\mathrm{d}\mathbf{v}) = M^T \bar{\psi}_\theta^{\mathbf{s}_{1,\mathbf{k}}}(\mathrm{d}\mathbf{v}) b_\theta(\mathbf{k}|\mathbf{s}_{1,\mathbf{k}}(\mathbf{v}))$$

*Proof of Corollary 2.* The corollary is a direct consequence of Lemmas 6, 7, and 8.  $\square$

**Corollary 3.** *For any  $\theta \in \Theta$ ,  $\mathbf{k} \in \llbracket M \rrbracket^T$  and  $\mathbf{v} \in \mathbb{Z}^{MT}$ ,*

$$(\pi_\theta \otimes \Phi_\theta)(\mathrm{d}\mathbf{v}) b_\theta(\mathbf{k}|\mathbf{v}) = (\pi_\theta \otimes \Phi_\theta)^{\mathbf{s}_{1,\mathbf{k}}}(\mathrm{d}\mathbf{v}) b_\theta(\mathbf{k}|\mathbf{s}_{1,\mathbf{k}}(\mathbf{v})).$$

*Proof of Corollary 3.* By Corollary 2, we have

$$(\pi_\theta \otimes \Phi_\theta)(\mathrm{d}\mathbf{v}) b_\theta(\mathbf{k}|\mathbf{v}) = M^T \bar{\psi}_\theta^{\mathbf{s}_{1,\mathbf{k}}}(\mathrm{d}\mathbf{v}) b_\theta(\mathbf{k}|\mathbf{s}_{1,\mathbf{k}}(\mathbf{v})) b_\theta(\mathbf{k}|\mathbf{v})$$

Also, note that, we have

$$\begin{aligned} \bar{\psi}_\theta^{\mathbf{s}_{1,\mathbf{k}}}(\mathrm{d}\mathbf{v}) b_\theta(\mathbf{k}|\mathbf{v}) &= \bar{\psi}_\theta^{\mathbf{s}_{1,\mathbf{k}}}(\mathrm{d}\mathbf{v}) b_\theta(\mathbf{1}|\mathbf{s}_{1,\mathbf{k}}(\mathbf{v})) \\ &= (\pi_\theta \otimes \Phi_\theta)^{\mathbf{s}_{1,\mathbf{k}}}(\mathrm{d}\mathbf{v}), \end{aligned}$$

where the second line follows from Lemma 6. Substituting the latter equation into the former, we conclude.  $\square$

In addition to the results above, the following lemma will be useful in Section C.1.

**Lemma 9.** *Let  $F : \mathbb{Z} \rightarrow \mathbb{R}$  be a real-valued function. Then, for any  $\theta \in \Theta$  we have*

$$\sum_{\mathbf{k} \in \llbracket M \rrbracket^T} \int_{\mathbb{Z}^{MT}} b_\theta(\mathbf{k}|\mathbf{v}) F(v^{(\mathbf{k})}) (\pi_\theta \otimes \Phi_\theta)(\mathrm{d}\mathbf{v}) = \sum_{\mathbf{k} \in \llbracket M \rrbracket^T} \int_{\mathbb{Z}^{MT}} F(v^{(\mathbf{k})}) b_\theta(\mathbf{k}|\mathbf{v}) \bar{\psi}_\theta(\mathrm{d}\mathbf{v}).$$

*Proof of Lemma 9.* First, notice that for any  $\mathbf{l} \in \llbracket M \rrbracket^T$ , we can write

$$\sum_{\mathbf{k} \in \llbracket M \rrbracket^T} b_\theta(\mathbf{k}|\mathbf{v}) F(v^{(\mathbf{k})}) = \sum_{\mathbf{k} \in \llbracket M \rrbracket^T} b_\theta(\mathbf{k}|\mathbf{s}_{1,\mathbf{l}}(\mathbf{v})) F(\mathbf{s}_{1,\mathbf{l}}(\mathbf{v})^{(\mathbf{k})}). \quad (49)$$

due to one-to-one correspondence of the paths in  $\mathbf{v}$  and  $\mathbf{s}_{1,\mathbf{l}}(\mathbf{v})$ . Combining this with Corollary 2 applied with path  $\mathbf{l}$ , we have

$$\begin{aligned} (\pi_\theta \otimes \Phi_\theta)(\mathrm{d}\mathbf{v}) \sum_{\mathbf{k} \in \llbracket M \rrbracket^T} F(v^{(\mathbf{k})}) b_\theta(\mathbf{k}|\mathbf{v}) \\ = M^T \bar{\psi}_\theta^{\mathbf{s}_{1,\mathbf{l}}}(\mathrm{d}\mathbf{v}) b_\theta(\mathbf{l}|\mathbf{s}_{1,\mathbf{l}}(\mathbf{v})) \sum_{\mathbf{k} \in \llbracket M \rrbracket^T} F(\mathbf{s}_{1,\mathbf{l}}(\mathbf{v})^{(\mathbf{k})}) b_\theta(\mathbf{k}|\mathbf{s}_{1,\mathbf{l}}(\mathbf{v})). \end{aligned}$$

Since the above holds for every  $\mathbf{l} \in \llbracket M \rrbracket^T$ , summing over  $\mathbf{l} \in \llbracket M \rrbracket^T$  and dividing by  $M^T$  results in

$$\begin{aligned} (\pi_\theta \otimes \Phi_\theta)(d\mathbf{v}) \sum_{\mathbf{k} \in \llbracket M \rrbracket^T} b_\theta(\mathbf{k}|\mathbf{v}) F(v^{(\mathbf{k})}) \\ = \sum_{\mathbf{l} \in \llbracket M \rrbracket^T} \bar{\psi}_\theta^{\mathbf{s}_{1,\mathbf{l}}}(\mathbf{v}) b_\theta(\mathbf{l}|\mathbf{s}_{1,\mathbf{l}}(\mathbf{v})) \sum_{\mathbf{k} \in \llbracket M \rrbracket^T} F(\mathbf{s}_{1,\mathbf{l}}(\mathbf{v})^{(\mathbf{k})}) b_\theta(\mathbf{k}|\mathbf{s}_{1,\mathbf{l}}(\mathbf{v})) \end{aligned}$$

Taking the integral of both sides over  $\mathbf{v}$ , we get

$$\begin{aligned} \int_{\mathbf{Z}^{MT}} (\pi_\theta \otimes \Phi_\theta)(d\mathbf{v}) \sum_{\mathbf{k} \in \llbracket M \rrbracket^T} F(v^{(\mathbf{k})}) b_\theta(\mathbf{k}|\mathbf{v}) \\ = \sum_{\mathbf{l} \in \llbracket M \rrbracket^T} \int_{\mathbf{Z}^{MT}} \bar{\psi}_\theta^{\mathbf{s}_{1,\mathbf{l}}}(\mathbf{v}) b_\theta(\mathbf{l}|\mathbf{s}_{1,\mathbf{l}}(\mathbf{v})) \sum_{\mathbf{k} \in \llbracket M \rrbracket^T} b_\theta(\mathbf{k}|\mathbf{s}_{1,\mathbf{l}}(\mathbf{v})) F(\mathbf{s}_{1,\mathbf{l}}(\mathbf{v})^{(\mathbf{k})}) \\ = \sum_{\mathbf{l} \in \llbracket M \rrbracket^T} \int_{\mathbf{Z}^{MT}} \bar{\psi}_\theta(\mathbf{v}) b_\theta(\mathbf{l}|\mathbf{v}) \sum_{\mathbf{k} \in \llbracket M \rrbracket^T} b_\theta(\mathbf{k}|\mathbf{v}) F(\mathbf{v}^{(\mathbf{k})}) \\ = \int_{\mathbf{Z}^{MT}} \bar{\psi}_\theta(\mathbf{v}) \sum_{\mathbf{k} \in \llbracket M \rrbracket^T} b_\theta(\mathbf{k}|\mathbf{v}) F(\mathbf{v}^{(\mathbf{k})}) \sum_{\mathbf{l} \in \llbracket M \rrbracket^T} b_\theta(\mathbf{l}|\mathbf{v}) \\ = \sum_{\mathbf{k} \in \llbracket M \rrbracket^T} \int_{\mathbf{Z}^{MT}} \bar{\psi}_\theta(\mathbf{v}) b_\theta(\mathbf{k}|\mathbf{v}) F(\mathbf{v}^{(\mathbf{k})}) \end{aligned}$$

where in the first and third lines we apply a change in the order of integration/summation, in the second line we apply a change of variables  $\mathbf{v} \rightarrow \mathbf{s}_{1,\mathbf{l}}(\mathbf{v})$ , whose Jacobian is 1, and the last line follows since  $b_\theta(\mathbf{l}|\mathbf{v})$  is a probability distribution for  $\mathbf{l}$ .  $\square$

## C.1 Unbiasedness for the acceptance ratio estimator of Algorithm 5

We provide a proof of Theorem 4 that states the unbiasedness for the acceptance ratio estimator of Algorithm 5.

*Proof of Theorem 4.* The expectation of  $r_{1,\mathbf{v}}(\theta, \vartheta; \zeta)$  with respect to the law of the mechanism described in Theorem 4 is

$$\begin{aligned} \int_{\mathbf{Z}^{MT}} \left[ \sum_{\mathbf{k} \in \llbracket M \rrbracket^T} r_{v^{(1)}, v^{(\mathbf{k})}}(\theta, \vartheta; \zeta) b_\zeta(\mathbf{k}|\mathbf{v}) \right] (\pi_\theta \otimes \Phi_\zeta)(d\mathbf{v}) \\ = \int_{\mathbf{Z}^{MT}} \left[ \sum_{\mathbf{k} \in \llbracket M \rrbracket^T} r_{v^{(1)}, v^{(\mathbf{k})}}(\theta, \vartheta; \zeta) b_\zeta(\mathbf{k}|\mathbf{v}) (\pi_\theta \otimes \Phi_\zeta)(d\mathbf{v}) \right] \end{aligned}$$

Let  $\gamma_\theta(z) := p_\theta(z, y)$  be the unnormalised density for  $\pi_\theta(z)$  so that  $\gamma_\theta(z) = \pi_\theta(z) \ell_\theta(y)$ .

Then, the term inside the sum on the RHS above can be written explicitly as

$$\begin{aligned}
r_{v^{(1)}, v^{(\mathbf{k})}}(\theta, \vartheta; \zeta) b_\zeta(\mathbf{k}|\mathbf{v})(\pi_\theta \otimes \Phi_\zeta)(d\mathbf{v}) &= \\
&= \frac{\eta(\vartheta)q(\vartheta, \theta)}{\eta(\theta)q(\theta, \vartheta)} \frac{\gamma_\zeta(v^{(1)})}{\gamma_\theta(v^{(1)})} \frac{\gamma_\vartheta(v^{(\mathbf{k})})}{\gamma_\zeta(v^{(\mathbf{k})})} b_\zeta(\mathbf{k}|\mathbf{v}) \frac{\pi_\theta(v^{(1)})}{\pi_\zeta(v^{(1)})} (\pi_\zeta \otimes \Phi_\zeta)(d\mathbf{v}) \\
&= \frac{\eta(\vartheta)q(\vartheta, \theta)}{\eta(\theta)q(\theta, \vartheta)} \frac{\ell_\zeta(y)}{\ell_\theta(y)} \frac{\gamma_\vartheta(v^{(\mathbf{k})})}{\gamma_\zeta(v^{(\mathbf{k})})} b_\zeta(\mathbf{k}|\mathbf{v})(\pi_\zeta \otimes \Phi_\zeta)(d\mathbf{v})
\end{aligned}$$

Integrating both sides over  $\mathbf{k}$  and  $\mathbf{v}$ , and applying Lemma 9, we have

$$\begin{aligned}
\int_{\mathbb{Z}^{MT}} \sum_{\mathbf{k} \in \llbracket M \rrbracket^T} r_{v^{(1)}, v^{(\mathbf{k})}}(\theta, \vartheta; \zeta) b_\zeta(\mathbf{k}|\mathbf{v})(\pi_\theta \otimes \Phi_\zeta)(d\mathbf{v}) \\
= \frac{\eta(\vartheta)q(\vartheta, \theta)}{\eta(\theta)q(\theta, \vartheta)} \frac{\ell_\zeta(y)}{\ell_\theta(y)} \sum_{\mathbf{k} \in \llbracket M \rrbracket^T} \int_{\mathbb{Z}^{MT}} \frac{\gamma_\vartheta(v^{(\mathbf{k})})}{\gamma_\zeta(v^{(\mathbf{k})})} b_\zeta(\mathbf{k}|\mathbf{v}) \psi_\zeta(d\mathbf{v}) \quad (50)
\end{aligned}$$

Next, we will show that, the integral on the right hand side in (50) is  $\frac{1}{M^T} \frac{\ell_\vartheta(y)}{\ell_\zeta(y)}$  for every  $\mathbf{k} \in \llbracket M \rrbracket^T$ . Indeed, fixing  $\mathbf{k} \in \llbracket M \rrbracket^T$ , we have, using the symmetry of  $\bar{\psi}_\zeta$ , Lemma 6, c.o.v  $\mathbf{v} \rightarrow \mathbf{s}_{1, \mathbf{k}}(\mathbf{v})$  and  $\mathbf{v} \sim (\pi_\zeta \otimes \Phi_\zeta)(\cdot) \Rightarrow v^{(1)} \sim \pi_\zeta(\cdot)$

$$\begin{aligned}
\int_{\mathbb{Z}^{MT}} \frac{\gamma_\vartheta(v^{(\mathbf{k})})}{\gamma_\zeta(v^{(\mathbf{k})})} b_\zeta(\mathbf{k}|\mathbf{v}) \bar{\psi}_\zeta(d\mathbf{v}) &= \int_{\mathbb{Z}^{MT}} \frac{\gamma_\vartheta(v^{(\mathbf{k})})}{\gamma_\zeta(v^{(\mathbf{k})})} b_\zeta(\mathbf{1}|\mathbf{s}_{1, \mathbf{k}}(\mathbf{v})) \bar{\psi}_\zeta^{\mathbf{s}_{1, \mathbf{k}}}(d\mathbf{v}) \\
&= \int_{\mathbb{Z}^{MT}} \frac{1}{M^T} \frac{\gamma_\vartheta(v^{(\mathbf{k})})}{\gamma_\zeta(v^{(\mathbf{k})})} (\pi_\zeta \otimes \Phi_\zeta)^{\mathbf{s}_{1, \mathbf{k}}}(d\mathbf{v}) \\
&= \int_{\mathbb{Z}^{MT}} \frac{1}{M^T} \frac{\gamma_\vartheta(v^{(1)})}{\gamma_\zeta(v^{(1)})} (\pi_\zeta \otimes \Phi_\zeta)(d\mathbf{v}) \\
&= \int_{\mathbb{Z}^{MT}} \frac{1}{M^T} \frac{\gamma_\vartheta(v^{(1)})}{\gamma_\zeta(v^{(1)})} \pi_\zeta(dv^{(1)}) \\
&= \frac{1}{M^T} \frac{\ell_\vartheta(y)}{\ell_\zeta(y)},
\end{aligned}$$

which does not depend on  $\mathbf{k}$ . Summing over  $M^T$  possible values of  $\mathbf{k}$ , and multiplying with the constant ratio on the right hand side of in (50), we obtain the expectation as

$$\begin{aligned}
\frac{\eta(\vartheta)q(\vartheta, \theta)}{\eta(\theta)q(\theta, \vartheta)} \frac{\ell_\zeta(y)}{\ell_\theta(y)} M^T \frac{1}{M^T} \frac{\ell_\vartheta(y)}{\ell_\zeta(y)} &= \frac{\eta(\vartheta)q(\vartheta, \theta)}{\eta(\theta)q(\theta, \vartheta)} \frac{\ell_\vartheta(y)}{\ell_\theta(y)} \\
&= r(\theta, \vartheta).
\end{aligned}$$

as required. □

## C.2 Proofs for Algorithm 5

### C.2.1 Acceptance ratio of Algorithm 5

The following lemmas can be verified by inspection.



**Lemma 10.** For any  $\theta, \vartheta, \zeta \in \Theta$ ,  $\mathbf{v} \in \mathbb{Z}^{MT}$ , and  $\mathbf{k} \in \llbracket M \rrbracket^T$ , we have  $r_{\mathbf{k}, \mathbf{v}}(\theta, \vartheta; \zeta) = r_{\mathbf{1}, \mathbf{s}_{\mathbf{1}, \mathbf{k}}(\mathbf{v})}(\theta, \vartheta; \zeta)$ .

**Lemma 11.** For any  $\mathbf{v} \in \mathbb{Z}^{MT}$ , and  $\mathbf{k}, \mathbf{l} \in \llbracket M \rrbracket^T$ , we have  $r_{\mathbf{k}, \mathbf{v}}(\theta, \vartheta; \theta) = r_{\mathbf{l}, \mathbf{v}}(\theta, \vartheta; \theta)$ .

We proceed to prove Theorem 5 that states the acceptance ratio of Algorithm 5.

*Proof of Theorem 5.* The joint distribution that corresponds to the moves of Algorithm 5 is

$$\begin{aligned} \hat{\pi}(d\xi) &= \frac{1}{2} \mathbb{I}_1(c) \pi(d\theta) q(\theta, d\vartheta) (\pi_\theta \otimes \Phi_{\zeta_1(\theta, \vartheta)})(d\mathbf{v}) \frac{r_{v(1), v(\mathbf{k})}(\theta, \vartheta; \zeta_1(\theta, \vartheta)) b_{\zeta_1(\theta, \vartheta)}(\mathbf{k}|\mathbf{v})}{r_{\mathbf{1}, \mathbf{v}}(\theta, \vartheta; \zeta_1(\theta, \vartheta))} \\ &\quad + \frac{1}{2} \mathbb{I}_2(c) \pi(d\theta) q(\theta, d\vartheta) (\pi_\theta \otimes \Phi_{\zeta_2(\theta, \vartheta)})(d\mathbf{v}) b_{\zeta_2(\theta, \vartheta)}(\mathbf{k}|\mathbf{v}). \end{aligned}$$

The proposed involution is  $\varphi(\theta, \vartheta, \mathbf{v}, \mathbf{k}, c) = (\vartheta, \theta, \mathbf{s}_{\mathbf{1}, \mathbf{k}}(\mathbf{v}), \mathbf{k}, 3 - c)$ . First, observe that, for any  $z, z' \in \mathbb{Z}^T$ , and  $\theta, \vartheta, \zeta \in \Theta$ , equation (34) can be rewritten as

$$\begin{aligned} \mathring{r}_{z, z'}(\theta, \vartheta; \zeta) &= \frac{q(\vartheta, \theta)}{q(\theta, \vartheta)} \frac{\eta(\vartheta)}{\eta(\theta)} \frac{\pi(d(\vartheta, z'))}{\pi(d(\zeta, z'))} \frac{\pi(d(\zeta, z))}{\pi(d(\theta, z))} \\ &= r(\theta, \vartheta) \frac{\pi_\vartheta(dz')}{\pi_\zeta(dz')} \frac{\pi_\zeta(dz)}{\pi_\theta(dz)}. \end{aligned} \quad (51)$$

where  $r(\theta, \vartheta)$  is defined in (32). Letting  $\zeta = \zeta_1(\theta, \vartheta) = \zeta_2(\vartheta, \theta)$ , when  $c = 1$ , we arrive at the acceptance ratio

$$\begin{aligned} \frac{\hat{\pi}^\varphi(d\xi)}{\hat{\pi}(d\xi)} &= r(\theta, \vartheta) \frac{(\pi_\vartheta \otimes \Phi_\zeta)^{\mathbf{s}_{\mathbf{1}, \mathbf{k}}}(\mathbf{v}) b_\zeta(\mathbf{k}|\mathbf{s}_{\mathbf{1}, \mathbf{k}}(\mathbf{v}))}{(\pi_\theta \otimes \Phi_\zeta)(d\mathbf{v}) \frac{b_\zeta(\mathbf{k}|\mathbf{v}) r_{v(1), v(\mathbf{k})}(\theta, \vartheta; \zeta)}{r_{\mathbf{1}, \mathbf{v}}(\theta, \vartheta; \zeta)}} \\ &= \frac{r(\theta, \vartheta)}{r_{v(1), v(\mathbf{k})}(\theta, \vartheta; \zeta)} \frac{(\pi_\vartheta \otimes \Phi_\zeta)^{\mathbf{s}_{\mathbf{1}, \mathbf{k}}}(\mathbf{v}) b_\zeta(\mathbf{k}|\mathbf{s}_{\mathbf{1}, \mathbf{k}}(\mathbf{v}))}{(\pi_\theta \otimes \Phi_\zeta)(d\mathbf{v}) b_\zeta(\mathbf{k}|\mathbf{v})} r_{\mathbf{1}, \mathbf{v}}(\theta, \vartheta; \zeta) \\ &= \frac{r(\theta, \vartheta)}{r_{v(1), v(\mathbf{k})}(\theta, \vartheta; \zeta)} \frac{\frac{\pi_\vartheta(dv(\mathbf{k}))}{\pi_\zeta(dv(\mathbf{k}))}}{\frac{\pi_\theta(dv(1))}{\pi_\zeta(dv(1))}} \frac{(\pi_\zeta \otimes \Phi_\zeta)^{\mathbf{s}_{\mathbf{1}, \mathbf{k}}}(\mathbf{v}) b_\zeta(\mathbf{k}|\mathbf{s}_{\mathbf{1}, \mathbf{k}}(\mathbf{v}))}{(\pi_\zeta \otimes \Phi_\zeta)(d\mathbf{v}) b_\zeta(\mathbf{k}|\mathbf{v})} r_{\mathbf{1}, \mathbf{v}}(\theta, \vartheta; \zeta) \\ &= r_{\mathbf{1}, \mathbf{v}}(\theta, \vartheta; \zeta), \end{aligned} \quad (52)$$

where for the last line we have used Corollary 2 and (51). For  $c = 2$ , upon using (7), we write

$$\begin{aligned} \mathring{r}(\theta, \vartheta, \mathbf{v}, \mathbf{k}, 2) &= \mathring{r}(\vartheta, \theta, \mathbf{s}_{\mathbf{1}, \mathbf{k}}(\mathbf{v}), \mathbf{k}, 1)^{-1} \\ &= r_{\mathbf{1}, \mathbf{s}_{\mathbf{1}, \mathbf{k}}(\mathbf{v})}(\vartheta, \theta; \zeta_1(\vartheta, \theta))^{-1} \\ &= r_{\mathbf{k}, \mathbf{v}}(\vartheta, \theta; \zeta_2(\theta, \vartheta))^{-1}, \end{aligned}$$

where the last line follows from Lemma 10, which concludes the proof.  $\square$

The analysis in the proof above not only bears an alternative proof of Theorem 4 on the unbiasedness of (35) but also implicitly proves Corollary 1; as we show below.

*Proof of Theorem 4.* By the equality of the first and last lines of (52), for any  $(\theta, \vartheta, \mathbf{v}, \mathbf{k}) \in \Theta^2 \times \mathbb{Z}^{MT} \times \llbracket M \rrbracket^T$ , we can write

$$(\pi_\theta \otimes \Phi_\zeta)(d\mathbf{v}) \frac{r_{v^{(1)}, v^{(k)}}(\theta, \vartheta; \zeta) b_\zeta(\mathbf{k}|\mathbf{v})}{\sum_{\mathbf{l} \in \llbracket M \rrbracket^T} r_{v^{(1)}, v^{(l)}}(\theta, \vartheta; \zeta) b_\zeta(\mathbf{l}|\mathbf{v})} \frac{r_{1, \mathbf{v}}(\theta, \vartheta; \zeta)}{r(\theta, \vartheta)} = (\pi_\vartheta \otimes \Phi_\zeta)^{\mathbf{s}_{1, \mathbf{k}}}(\mathbf{v}) b_\zeta(\mathbf{k}|\mathbf{s}_{1, \mathbf{k}}(\mathbf{v})). \quad (53)$$

Integrating both sides with respect to all the variables except  $\theta$  and  $\vartheta$  leads to

$$\int_{\mathbb{Z}^{MT}} r_{1, \mathbf{v}}(\theta, \vartheta; \zeta) (\pi_\theta \otimes \Phi_\zeta)(d\mathbf{v}) = r(\theta, \vartheta)$$

upon noticing that  $r_{1, \mathbf{v}}(\theta, \vartheta; \zeta)$  does not depend on  $\mathbf{k}$  and the right hand side is a probability distribution for  $(\mathbf{v}, \mathbf{k})$ . Recalling  $z = v^{(1)}$  and noting that  $(\pi_\theta \otimes \Phi_\zeta)(d\mathbf{v})$  is exactly the distribution of the mechanism described in Theorem 4 that generates  $r_{1, \mathbf{v}}(\theta, \vartheta; \zeta)$ , we prove Theorem 4.  $\square$

*Proof of Corollary 1.* Similarly to the previous proof, we make use of (52). However, this time we write the identity in (53) for  $(\vartheta, \theta, \mathbf{s}_{1, \mathbf{k}}(\mathbf{v}), \mathbf{k})$  to obtain

$$(\pi_\theta \otimes \Phi_\zeta)(d\mathbf{v}) b_\zeta(\mathbf{k}|\mathbf{v}) \frac{r(\vartheta, \theta)}{r_{\mathbf{k}, \mathbf{v}}(\vartheta, \theta; \zeta)} = (\pi_\vartheta \otimes \Phi_\zeta)^{\mathbf{s}_{1, \mathbf{k}}}(\mathbf{v}) \frac{r_{v^{(k)}, v^{(1)}}(\vartheta, \theta; \zeta) b_\zeta(\mathbf{k}|\mathbf{s}_{1, \mathbf{k}}(\mathbf{v}))}{\sum_{\mathbf{l} \in \llbracket M \rrbracket^T} r_{v^{(k)}, \mathbf{s}_{1, \mathbf{k}}(\mathbf{v})^{(l)}}(\vartheta, \theta; \zeta) b_\zeta(\mathbf{l}|\mathbf{s}_{1, \mathbf{k}}(\mathbf{v}))}$$

Again, integrating both sides with respect to all the variables  $\mathbf{v}, \mathbf{k}$ , we get

$$\int_{\mathbb{Z}^{MT}} \sum_{\mathbf{k} \in \llbracket M \rrbracket^T} (1/r_{\mathbf{k}, \mathbf{v}}(\vartheta, \theta; \zeta)) b_\zeta(\mathbf{k}|\mathbf{v}) (\pi_\theta \otimes \Phi_\zeta)(d\mathbf{v}) = r(\theta, \vartheta).$$

Since  $1/r_{\mathbf{k}, \mathbf{v}}(\vartheta, \theta; \zeta)$  is the estimator in question in Corollary 1 and  $(\pi_\theta \otimes \Phi_\zeta)(d\mathbf{v}) b_\zeta(\mathbf{k}|\mathbf{v})$  is exactly the distribution of the described mechanism that generates it, we prove Corollary 1.  $\square$

### C.2.2 Delayed rejection step for Algorithm 5

When the delayed rejection step is included in Algorithm 5, the algorithm targets the modified joint distribution for  $\check{\xi} = (\xi, \mathbf{l})$ , defined as

$$\tilde{\pi}(d\check{\xi}) = \tilde{\pi}(d\xi) [\mathbb{I}_1(c) b_\theta(\mathbf{l}|\mathbf{v}) + \mathbb{I}_2(c) b_\vartheta(\mathbf{l}|\mathbf{s}_{1, \mathbf{k}}(\mathbf{v}))].$$

The conditional probability of the extra variable  $\mathbf{l}$  is simply the backward sampling probability of the cSMC kernel run at  $\theta$ . Now one iteration of the algorithm can be thought of as a two-stage procedure, where the first stage is the regular MHAAR update and the second stage is executed conditional on the result of the former. The two moves are given below.

1. In the first stage, MHAAR attempts a transition for the joint variable  $\check{\xi} = (\xi, \mathbf{l})$  as

$$\check{\varphi}_1(\xi, \mathbf{l}) := (\varphi(\xi), \mathbf{l}).$$

As  $\tilde{\varphi}_1$  is an involution, it yields the acceptance ratio as

$$\begin{aligned}\tilde{r}_1(\check{\xi}) &:= \frac{\tilde{\pi}^{\tilde{\varphi}_1}(\mathrm{d}\check{\xi})}{\tilde{\pi}(\mathrm{d}\check{\xi})} \\ &= \frac{\overset{\circ}{\pi}^{\varphi}(\mathrm{d}\xi)}{\overset{\circ}{\pi}(\mathrm{d}\xi)} \frac{b_{\theta}(\mathbf{l}|\mathfrak{s}_{1,\mathbf{k}} \circ \mathfrak{s}_{1,\mathbf{k}}(\mathbf{v}))}{b_{\theta}(\mathbf{l}|\mathbf{v})} \\ &= \frac{\overset{\circ}{\pi}^{\varphi}(\mathrm{d}\xi)}{\overset{\circ}{\pi}(\mathrm{d}\xi)} = \overset{\circ}{r}(\xi)\end{aligned}$$

which is exactly the same acceptance ratio we would have for the basic version of the algorithm that does not have the delayed rejection step. As it can be seen from the above derivation,  $\mathbf{l}$  does not need to be sampled at this stage, i.e., prior to the delayed rejection step, since the acceptance probability is independent of  $\mathbf{l}$ . The delayed rejection step can be performed by the following involution.

2. The proposed involution of delayed rejection is

$$\check{\varphi}_2(\xi, \mathbf{l}) := \begin{cases} (\theta, \vartheta, \mathfrak{s}_{1,\mathbf{l}}(\mathbf{v}), \mathfrak{r}_1(\mathbf{k}), c, \mathbf{l}), & c = 1, \\ (\xi, \mathbf{l}), & c = 2, \end{cases}$$

where  $\mathfrak{r}_1(\mathbf{k})$  is defined in equation (43). That is, we only perform the delayed rejection move when  $c = 1$  and  $\zeta_1(\theta, \vartheta) = \theta$  is chosen for the intermediate distribution. The acceptance ratio of this move can be written as

$$\tilde{r}_2(\check{\xi}) = \frac{\tilde{\pi}^{\check{\varphi}_2}(\mathrm{d}\check{\xi})}{\tilde{\pi}(\mathrm{d}\check{\xi})} \frac{1 - \min\{1, \tilde{r}_1 \circ \check{\varphi}_2(\check{\xi})\}}{1 - \min\{1, \tilde{r}_1(\check{\xi})\}} \quad (54)$$

**Theorem 7.** Assume  $\zeta_1(\theta, \vartheta) = \theta$ . Then,  $\tilde{r}_2(\check{\xi}) = 1$ .

*Proof of Theorem 7.* We prove the theorem by showing that both ratios in (54) are equal to 1 if  $\zeta_1(\theta, \vartheta) = \theta$ . Assume that  $\zeta_1(\theta, \vartheta) = \theta$ . When  $c = 1$ ,

$$\begin{aligned}\frac{\tilde{\pi}^{\check{\varphi}_2}(\mathrm{d}\check{\xi})}{\tilde{\pi}(\mathrm{d}\check{\xi})} &= r(\theta, \vartheta) \frac{(\pi_{\theta} \otimes \Phi_{\theta})^{\mathfrak{s}_{1,\mathbf{l}}}(\mathrm{d}\mathbf{v}) b_{\theta}(\mathfrak{r}_1(\mathbf{k})|\mathfrak{s}_{1,\mathbf{l}}(\mathbf{v})) b_{\theta}(\mathbf{l}|\mathfrak{s}_{1,\mathbf{l}}(\mathbf{v}))}{(\pi_{\theta} \otimes \Phi_{\theta})(\mathrm{d}\mathbf{v}) b_{\theta}(\mathbf{k}|\mathbf{v}) b_{\theta}(\mathbf{l}|\mathbf{v})} \\ &= r(\theta, \vartheta) \frac{(\pi_{\theta} \otimes \Phi_{\theta})^{\mathfrak{s}_{1,\mathbf{l}}}(\mathrm{d}\mathbf{v}) b_{\theta}(\mathbf{l}|\mathfrak{s}_{1,\mathbf{l}}(\mathbf{v}))}{(\pi_{\theta} \otimes \Phi_{\theta})(\mathrm{d}\mathbf{v}) b_{\theta}(\mathbf{l}|\mathbf{v})} \frac{b_{\theta}(\mathfrak{r}_1(\mathbf{k})|\mathfrak{s}_{1,\mathbf{l}}(\mathbf{v}))}{b_{\theta}(\mathbf{k}|\mathbf{v})},\end{aligned}$$

and all of the ratios are equal to 1. Moreover, the ratio involving the rejection probabilities is

$$\begin{aligned}\frac{1 - \min\{1, \tilde{r}_1 \circ \check{\varphi}_2(\xi, \mathbf{l})\}}{1 - \min\{1, \tilde{r}_1(\xi, \mathbf{l})\}} &= \frac{1 - \min\{1, \overset{\circ}{r}(\theta, \vartheta, \mathfrak{s}_{1,\mathbf{l}}(\mathbf{v}), \mathfrak{r}_1(\mathbf{k}), 1)\}}{1 - \min\{1, \overset{\circ}{r}(\theta, \vartheta, \mathbf{v}, \mathbf{k}, 1)\}} \\ &= \frac{1 - \min\{1, r_{1,\mathbf{v}}(\theta, \vartheta)\}}{1 - \min\{1, r_{1,\mathbf{v}}(\theta, \vartheta)\}} \\ &= 1,\end{aligned}$$

where the second line is by Lemma 10, and the last line is by Lemma 11. When  $c = 2$ , the move  $\check{\varphi}_2$  imposes no change, so the acceptance ratio is trivially equal to 1.  $\square$

Note that the conditions  $c = 1$  and  $\zeta_1(\theta, \vartheta) = \theta$  are critical here: The proposed update of delayed rejection does not change the sample for  $\theta$  but changes the sample for  $z$  via backward sampling at  $\theta$  conditional on the particles generated by an cSMC kernel run at  $\zeta_1(\theta, \vartheta)$ . Accepting this proposal with probability 1 preserves invariance only if  $\zeta_1(\theta, \vartheta) = \theta$ . We could, in theory, have a similar delayed rejection step when  $c = 2$  if  $\zeta_1(\theta, \vartheta) = \vartheta$ . However, the choice  $\zeta_1(\theta, \vartheta) = \vartheta$  is senseless because it disables all the averaging in the MHAAR algorithm, see equations (34) and (35).

### C.3 The subsampled version of MHAAR-RB for SSM

The subsampled version of MHAAR-RB-SSM, named MHAAR-S-SSM, which was mentioned in Section 4.2.3 is presented in Algorithm 7. Like in MHAAR-RB-SSM, refreshing  $z$  is also possible in Algorithm 7 as well, but in a different fashion, see the step labeled ‘optional’. Specifically, when  $c = 1$  and  $\zeta_1(\theta, \vartheta) = \theta$ , one can randomly swap  $z$  with  $u^{(i)}$  with a probability  $1/N$  for all  $i = 1, \dots, N$ , owing to exchangeability arguments. Note that this is not a delayed rejection step and the swapping has to be performed before making a decision, as it affects the acceptance ratio. However the computational cost of swapping two paths is negligible. We explain why this move preserves invariance in Appendix C.3.2.

---

**Algorithm 7:** MHAAR-S for SSM - reduced computation via subsampling

---

**Input:** Current sample  $(\theta, z)$   
**Output:** New sample

- 1 Sample  $\vartheta \sim q(\theta, \cdot)$  and  $c \sim \text{Unif}(\{1, 2\})$ , and set  $\zeta = \zeta_c(\theta, \vartheta)$ .
- 2 **if**  $c = 1$  **then**
- 3     Run a cSMC( $M, \zeta, z$ ) to obtain the particles  $\mathbf{v}$ .
- 4     Sample  $u^{(1)}, \dots, u^{(N)} \stackrel{\text{iid}}{\sim} \sum_{\mathbf{l} \in \llbracket M \rrbracket^T} \Phi_{\zeta}(\mathbf{l} | \mathbf{v}) \delta_{v^{(1)}}(\cdot)$ .
- 5     **if**  $\zeta = \theta$  **then**
- 6         Swap  $z$  with  $u^{(j)}$  where  $j \sim \text{Unif}(\llbracket N \rrbracket)$ . (optional refreshment of  $z$ )
- 7     Sample  $k \sim \mathcal{P}(r_{z, u^{(1)}}(\theta, \vartheta; \zeta), \dots, r_{z, u^{(N)}}(\theta, \vartheta; \zeta))$  and set  $z' = u^{(k)}$ .
- 8     Return  $(\vartheta, z')$  with probability  $\min\{1, r_{z, u}^N(\theta, \vartheta; \zeta)\}$ ; otherwise return  $(\theta, z)$ .
- 9 **else**
- 10     Run a cSMC( $M, \zeta, z$ ) to obtain particles  $\mathbf{v}$ .
- 11     Sample  $u^{(1)}, \dots, u^{(N)} \stackrel{\text{iid}}{\sim} \sum_{\mathbf{l} \in \llbracket M \rrbracket^T} \Phi_{\zeta}(\mathbf{l} | \mathbf{v}) \delta_{v^{(1)}}(\cdot)$ .
- 12     Sample  $k \sim \text{Unif}(\llbracket N \rrbracket)$ , set  $z' = u^{(k)}$ , and change  $u^{(k)} = z$ .
- 13     Return  $(\vartheta, z')$  with probability  $\min\{1, 1/r_{z', u}^N(\vartheta, \theta; \zeta)\}$ ; otherwise return  $(\theta, z)$ .

---

#### C.3.1 Reversibility of Algorithm 7

Next, we show the reversibility of Algorithm 7 that uses a subsampled version of the Rao-Blackwellised acceptance ratio estimator.

For any  $\theta \in \Theta$ , suppose  $u^{(0)} \sim \pi_{\theta}(\cdot)$  and let  $u^{(1)}, \dots, u^{(N)}$  be  $N$  paths drawn via backward sampling following cSMC at  $\zeta$  conditioned on  $u^{(0)}$ . Then the joint distribution

of  $\mathbf{u} := (u^{(0)}, \dots, u^{(N)})$  can be written as

$$R_{\theta, \zeta}(d\mathbf{u}) = \int_{\mathbf{Z}^{MT}} \left\{ [(\pi_\theta \otimes \Phi_\zeta)(d\mathbf{v}) \delta_{v(1)}(du^{(0)})] \prod_{i=1}^N \left[ \sum_{\mathbf{k} \in \llbracket M \rrbracket^T} b_\zeta(\mathbf{k}|\mathbf{v}) \delta_{v(\mathbf{k})}(du^{(i)}) \right] \right\}.$$

**Lemma 12.** *The following hold for  $R_{\theta, \zeta}(d(u^{(0)}, \dots, u^{(N)}))$ :*

1. *The marginal distribution of  $u^{(0)}$  is  $\pi_\theta(\cdot)$ .*
2. *When  $\theta = \zeta$ , the variables  $u^{(0)}, u^{(1)}, \dots, u^{(N)}$  are exchangeable and share  $\pi_\theta(\cdot)$  as their marginal distribution.*
3.  *$R_{\theta, \zeta}(d\mathbf{u}) = \frac{\pi_\theta(du^{(0)})}{\pi_\zeta(du^{(0)})} R_{\zeta, \zeta}(d\mathbf{u})$ .*

*Proof of Lemma 12.* The claims in the lemma can be proven by considering the joint distribution

$$\bar{R}_{\theta, \zeta}(d(\mathbf{u}, \mathbf{v}, \mathbf{k}_0, \dots, \mathbf{k}_N)) := \frac{\pi_\theta(du^{(0)})}{\pi_\zeta(du^{(0)})} \bar{\psi}_\zeta(d\mathbf{v}) \prod_{i=0}^N b_\zeta(\mathbf{k}_i|\mathbf{v}) \delta_{v(\mathbf{k}_i)}(du^{(i)}). \quad (55)$$

First, we show that  $\bar{R}_{\theta, \zeta}(d(\mathbf{u}, \mathbf{v}, \mathbf{k}_0, \dots, \mathbf{k}_N))$  is a distribution whose marginal distribution for  $\mathbf{u}$  is  $R_{\theta, \zeta}(\mathbf{u})$ . For this, first note the identity

$$\frac{\pi_\theta(du^{(0)})}{\pi_\zeta(du^{(0)})} \bar{\psi}_\zeta(d\mathbf{v}) b_\zeta(\mathbf{k}_0|\mathbf{v}) \delta_{v(\mathbf{k}_0)}(du^{(0)}) = \frac{\pi_\theta(du^{(0)})}{\pi_\zeta(du^{(0)})} \bar{\psi}_\zeta^{\mathbf{s}_{\mathbf{1}, \mathbf{k}_0}}(d\mathbf{v}) b_\zeta(\mathbf{1}|\mathbf{s}_{\mathbf{1}, \mathbf{k}_0}(\mathbf{v})) \delta_{v(\mathbf{1})}(du^{(0)}), \quad (56)$$

which follows from Lemmas 7 and 8. Next, integrating the RHS of (56) with respect to  $\mathbf{k}_0$  and  $\mathbf{v}$ , we obtain

$$\begin{aligned} & \frac{\pi_\theta(du^{(0)})}{\pi_\zeta(du^{(0)})} \sum_{\mathbf{k}_0 \in \llbracket M \rrbracket^T} \int_{\mathbf{Z}^{MT}} \bar{\psi}_\zeta^{\mathbf{s}_{\mathbf{1}, \mathbf{k}_0}}(d\mathbf{v}) b_\zeta(\mathbf{1}|\mathbf{s}_{\mathbf{1}, \mathbf{k}_0}(\mathbf{v})) \delta_{v(\mathbf{1})}(du^{(0)}) \\ &= \frac{\pi_\theta(du^{(0)})}{\pi_\zeta(du^{(0)})} \int_{\mathbf{Z}^{MT}} M^T \bar{\psi}_\zeta(d\mathbf{v}) b_\zeta(\mathbf{1}|\mathbf{v}) \delta_{v(\mathbf{1})}(du^{(0)}) \\ &= \frac{\pi_\theta(du^{(0)})}{\pi_\zeta(du^{(0)})} \int_{\mathbf{Z}^{MT}} (\pi_\zeta \otimes \Phi_\zeta)(d\mathbf{v}) \delta_{v(1)}(du^{(0)}) \\ &= \int_{\mathbf{Z}^{MT}} (\pi_\theta \otimes \Phi_\zeta)(d\mathbf{v}) \delta_{v(1)}(du^{(0)}) \quad (57) \\ &= \pi_\theta(du^{(0)}), \quad (58) \end{aligned}$$

where in the second line we use a change of variable  $\mathbf{v} \rightarrow \mathbf{s}_{\mathbf{1}, \mathbf{k}_0}(\mathbf{v})$  and end up with the same expression for all  $\mathbf{k}_0$ , and the third line is by Lemma 6. Using (57) together with

(56), we have

$$\begin{aligned}
& \int_{\mathbf{Z}^{MT}} \sum_{\mathbf{k}_{0:N} \in \llbracket M \rrbracket^{TN}} \bar{R}_{\theta, \zeta}(\mathrm{d}(\mathbf{u}, \mathbf{v}, \mathbf{k}_0, \dots, \mathbf{k}_N)) \\
&= \int_{\mathbf{Z}^{MT}} (\pi_\theta \otimes \Phi_\zeta)(\mathrm{d}\mathbf{v}) \delta_{v^{(1)}}(\mathrm{d}u^{(0)}) \sum_{\mathbf{k}_{1:N} \in \llbracket M \rrbracket^{TN}} \prod_{i=1}^N [b_\zeta(\mathbf{k}_i | \mathbf{v}) \delta_{v^{(k_i)}}(\mathrm{d}u^{(i)})] \\
&= \int_{\mathbf{Z}^{MT}} \left\{ [(\pi_\theta \otimes \Phi_\zeta)(\mathrm{d}\mathbf{v}) \delta_{v^{(1)}}(\mathrm{d}u^{(0)})] \prod_{i=1}^N \left[ \sum_{\mathbf{k} \in \llbracket M \rrbracket^T} b_\zeta(\mathbf{k} | \mathbf{v}) \delta_{v^{(k)}}(\mathrm{d}u^{(i)}) \right] \right\} \\
&= R_{\theta, \zeta}(\mathrm{d}\mathbf{u}).
\end{aligned}$$

Now, we can proceed to proving the claims in the lemma. The first claim can be proven by integrating (55) with respect to  $\mathbf{k}_1, \dots, \mathbf{k}_N, u^{(1)}, \dots, u^{(N)}$  and then with respect to  $\mathbf{k}_0$  and  $\mathbf{v}$ , where in the latter step we use (58). For the second claim, observe that when  $\theta = \zeta$  we have

$$\bar{R}_{\theta, \theta}(\mathrm{d}(\mathbf{u}, \mathbf{v}, \mathbf{k}_0, \dots, \mathbf{k}_N)) := \bar{\psi}_\theta(\mathrm{d}\mathbf{v}) \prod_{i=0}^N b_\theta(\mathbf{k}_i | \mathbf{v}) \delta_{v^{(k_i)}}(\mathrm{d}u^{(i)}).$$

Taking the integral of both sides with respect to  $\mathbf{v}$  and  $\mathbf{k}_0, \dots, \mathbf{k}_N$ , we have

$$R_{\theta, \theta}(\mathrm{d}\mathbf{u}) = \int_{\mathbf{Z}^{MT}} \bar{\psi}_\theta(\mathrm{d}\mathbf{v}) \prod_{i=0}^N \sum_{\mathbf{k} \in \llbracket M \rrbracket^T} b_\theta(\mathbf{k} | \mathbf{v}) \delta_{v^{(k)}}(\mathrm{d}u^{(i)}) \quad (59)$$

and the exchangeability of  $u^{(0)}, u^{(1)}, \dots, u^{(N)}$  is obvious from the symmetry in (59). Moreover, due to exchangeability, since  $u^{(0)}$  has marginal  $\pi_\theta(\mathrm{d}u^{(0)})$ , so do  $u^{(1)}, \dots, u^{(N)}$ . For the third claim, note the relation

$$\bar{R}_{\theta, \zeta}(\mathrm{d}(\mathbf{u}, \mathbf{v}, \mathbf{k}_0, \dots, \mathbf{k}_N)) := \frac{\pi_\theta(\mathrm{d}u^{(0)})}{\pi_\zeta(\mathrm{d}u^{(0)})} \bar{R}_{\zeta, \zeta}(\mathrm{d}(\mathbf{u}, \mathbf{v}, \mathbf{k}_0, \dots, \mathbf{k}_N))$$

from (55). Taking the integral of both sides over  $\mathbf{v}, \mathbf{k}_0, \dots, \mathbf{k}_N$ , we have the claimed equality.  $\square$

**Theorem 8.** *The transition probability of Algorithm 7 satisfies detailed balance with respect to  $\pi(\mathrm{d}(\theta, z))$ .*

*Proof of Theorem 8.* The joint distribution corresponding to the moves of Algorithm 7 can be shown to target the joint distribution for  $\xi := (\theta, \vartheta, \mathbf{u}, k, c)$ , defined as

$$\begin{aligned}
\dot{\pi}(\mathrm{d}\xi) &:= \frac{1}{2} \mathbb{I}_1(c) \pi(\mathrm{d}\theta) q(\theta, \mathrm{d}\vartheta) R_{\theta, \zeta_1(\theta, \vartheta)}(\mathrm{d}\mathbf{u}) \frac{r_{u^{(0)}, u^{(k)}}(\theta, \vartheta; \zeta_1(\theta, \vartheta))}{\sum_{i=1}^N r_{u^{(0)}, u^{(i)}}(\theta, \vartheta; \zeta_1(\theta, \vartheta))} \\
&\quad + \frac{1}{2} \mathbb{I}_2(c) \pi(\mathrm{d}\theta) q(\theta, \mathrm{d}\vartheta) R_{\theta, \zeta_2(\theta, \vartheta)}(\mathrm{d}\mathbf{u}) \frac{1}{N}.
\end{aligned}$$

where the latent variable is embedded in  $\mathbf{u}$  as  $z = u^{(0)}$ . Then, Lemma 12, the marginal for  $(\theta, u^{(0)})$  is  $\pi(x)$ . The proposed involution is

$$\varphi(\theta, \vartheta, \mathbf{u}, k, c) := (\vartheta, \theta, \mathfrak{s}_{0,k}(\mathbf{u}), k, 3 - c),$$

where  $\mathfrak{s}_{0,k}(\mathbf{u})$  is an operator that swaps  $u^{(0)}$  and  $u^{(k)}$  in  $\mathbf{u}$ . Next, we derive the acceptance ratios

$$\mathring{r}(\xi) := \frac{\mathring{\pi}^\varphi(d\xi)}{\mathring{\pi}(d\xi)}$$

for  $c = 1$  and  $c = 2$ . When  $c = 1$ , we have

$$\begin{aligned} \mathring{\pi}(d\xi) &= \pi(\theta)q(\theta, d\vartheta) \frac{\pi_\theta(du^{(0)})}{\pi_{\zeta_1(\theta, \vartheta)}(du^{(0)})} R_{\zeta_1(\theta, \vartheta), \zeta_1(\theta, \vartheta)}(d\mathbf{u}) \frac{r_{u^{(0)}, u^{(k)}}(\theta, \vartheta; \zeta_1(\theta, \vartheta))}{\sum_{i=1}^N r_{u^{(0)}, u^{(i)}}(\theta, \vartheta; \zeta_1(\theta, \vartheta))}, \\ \mathring{\pi}^\varphi(d\xi) &= \frac{1}{2} \pi(d\vartheta)q(\vartheta, d\theta) \frac{\pi_\vartheta(du^{(k)})}{\pi_{\zeta_1(\theta, \vartheta)}(du^{(k)})} R_{\zeta_1(\theta, \vartheta), \zeta_1(\theta, \vartheta)}^{\mathfrak{s}_{0,k}}(d\mathbf{u}) \frac{1}{N} \\ &= \frac{1}{2} \pi(d\vartheta)q(\vartheta, d\theta) \frac{\pi_\vartheta(du^{(k)})}{\pi_{\zeta_1(\theta, \vartheta)}(du^{(k)})} R_{\zeta_1(\theta, \vartheta), \zeta_1(\theta, \vartheta)}(d\mathbf{u}) \frac{1}{N}, \end{aligned}$$

where we have used Lemma 12 in the lines of both equations. Noting (51), we conclude that, for  $c = 1$ ,

$$\mathring{r}(\xi) = \frac{\mathring{\pi}^\varphi(d\xi)}{\mathring{\pi}(d\xi)} = \frac{1}{N} \sum_{i=1}^N r_{u^{(0)}, u^{(i)}}(\theta, \vartheta; \zeta_1(\theta, \vartheta)).$$

For  $c = 2$ , we use (7) to get

$$\mathring{r}(\theta, \vartheta, \mathbf{u}, k, 2) = \left[ \frac{1}{N} \sum_{i=0, i \neq k}^N r_{u^{(k)}, u^{(i)}}(\vartheta, \theta; \zeta_2(\theta, \vartheta)) \right]^{-1}$$

□

### C.3.2 Refreshing the latent variable in Algorithm 7

As MHAAR-S-SSM in Algorithm 7 suggests, we consider refreshing  $z$  only when  $c = 1$  and  $\zeta_1(\theta, \vartheta) = \theta$ . When  $c = 1$ , one iteration of the modified algorithm can be stated as follows: Given  $x = (\theta, z)$ ,

1. Sample  $c \sim \text{Unif}(\{1, 2\})$ , set  $u^{(0)} = z$  and sample  $N$  paths  $(u^{(1)}, \dots, u^{(N)})$  using a single cSMC conditioned on  $u^{(0)}$ .
2. If  $c = 1$ , perform a random swap  $u^{(0)} \leftrightarrow u^{(i)}$  with probability  $1/N$  for all  $i = 1, \dots, N$ .
3. Sample  $k$  with probability proportional to  $r_{u^{(0)}, u^{(k)}}(\theta, \vartheta; \theta)$ .
4. Propose and accept/reject the move  $(\theta, \vartheta, \mathbf{u}, k, 1) \rightarrow (\vartheta, \theta, \mathfrak{s}_{0,k}(\mathbf{u}), k, 2)$ .



(For practical reasons, the order of steps 3 and 4 can be reversed.) The step that refreshes  $z$  is the second step. By the exchangeability result for  $R_{\theta,\theta}$  in Lemma 12, step 2 can be shown to target the conditional distribution (with respect to  $\overset{\circ}{\pi}$ ) of  $\mathbf{u}$  given  $\theta, \vartheta$ , and  $c$ , while  $\mathbf{k}$  is marginalised out. Therefore, the fact that this swap move preserves invariance of  $\overset{\circ}{\pi}$  follows from similar arguments for a collapsed Gibbs move.

Note that step 2 is not a delayed rejection step and it needs to be implemented before steps 3 and 4. However, this is not an issue computationally, since the computational complexity of the step is  $\mathcal{O}(1)$ .